

Как мы Change Data Capture делали

Александр Деулин
Василий Тюбек



HighLoad++
Весна 2021



Tarantool на Highload



Деулин Александр

aleksandr.deulin@nexign.com



Тюбек Василий

v.tyubek@corp.mail.ru

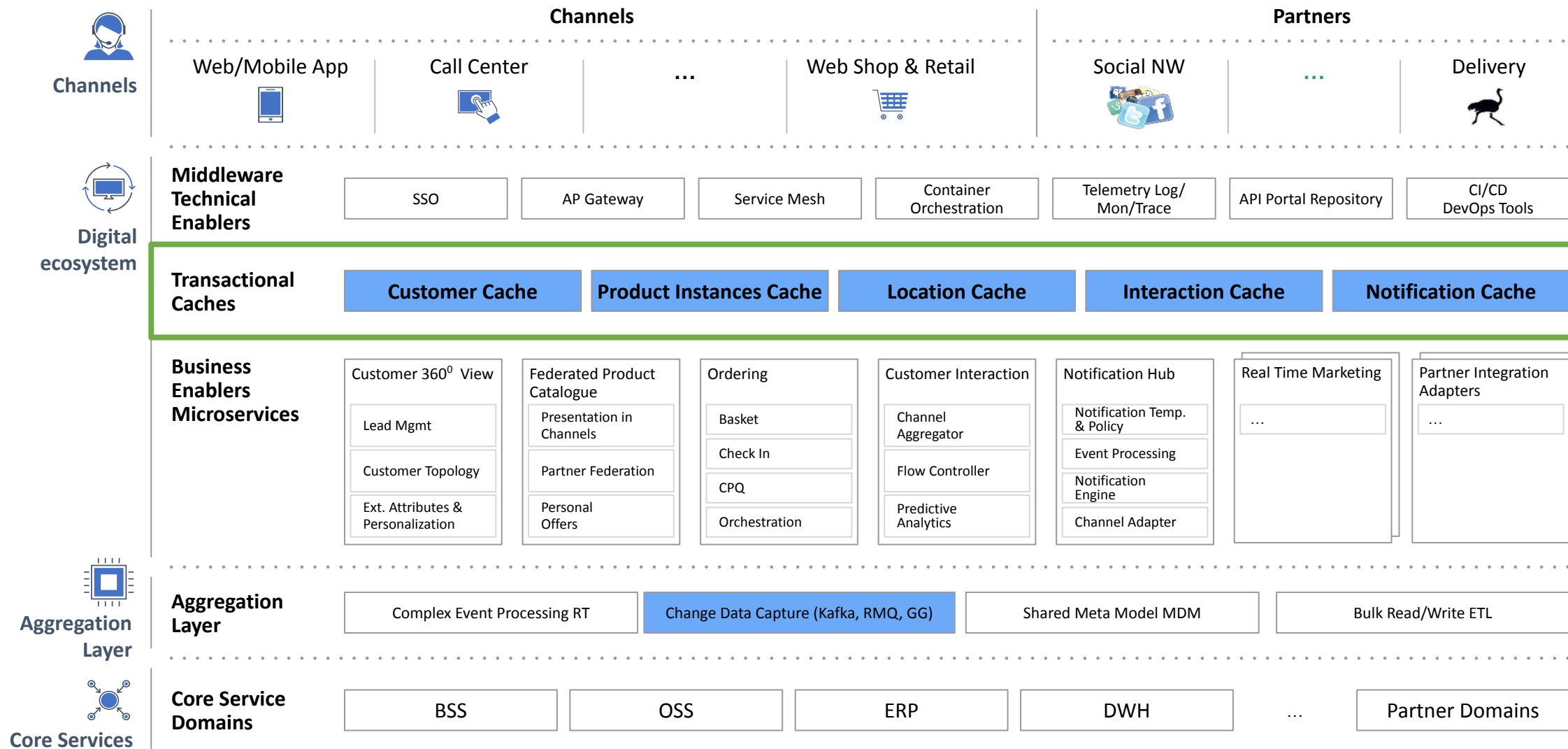
В предыдущем сезоне “Микросервисов”...



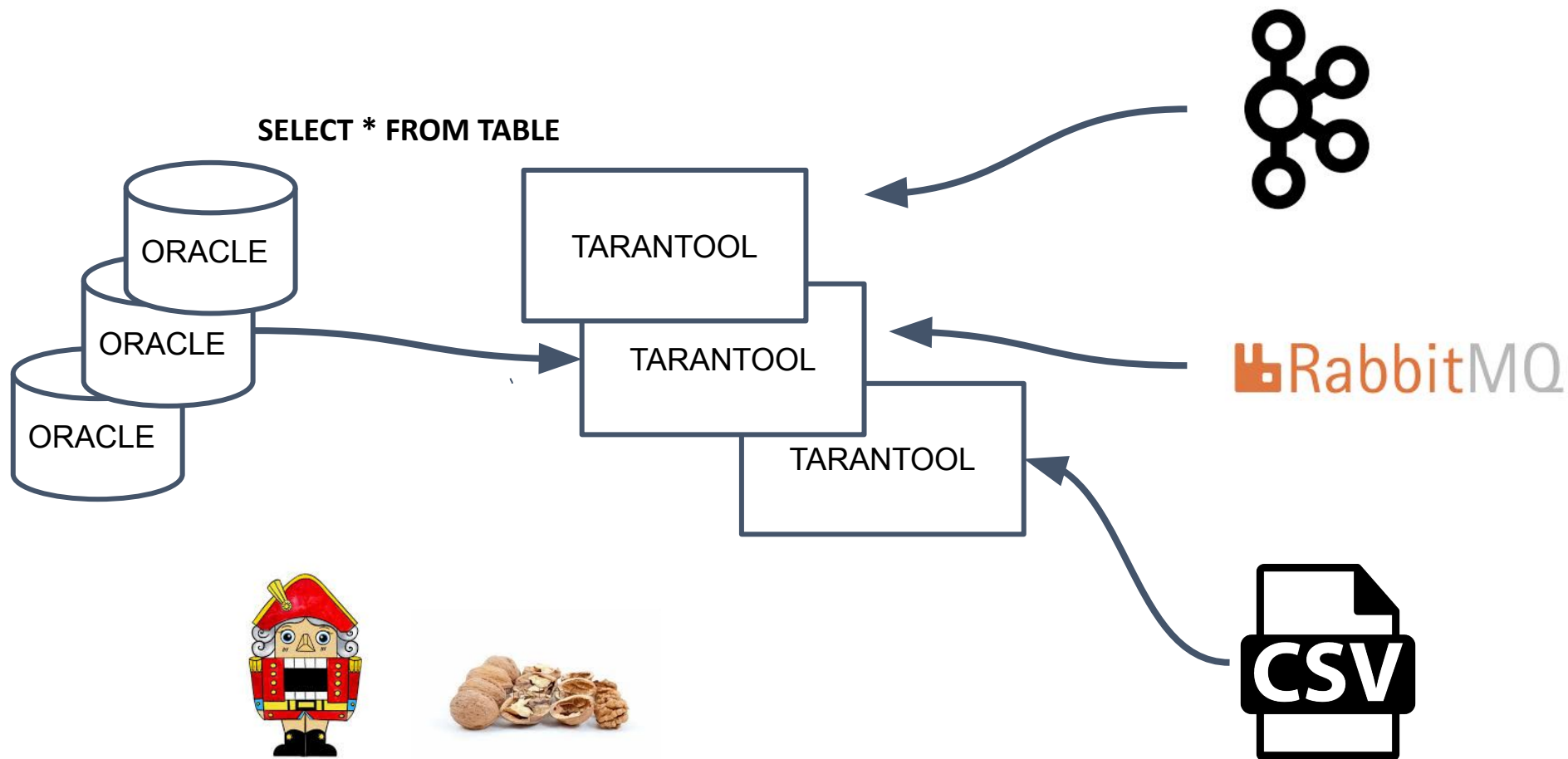
Применение микросервисов в высоконагруженном биллинге

<https://www.highload.ru/siberia/2019/abstracts/5384>

Все началось с кэшей



Вначале было просто

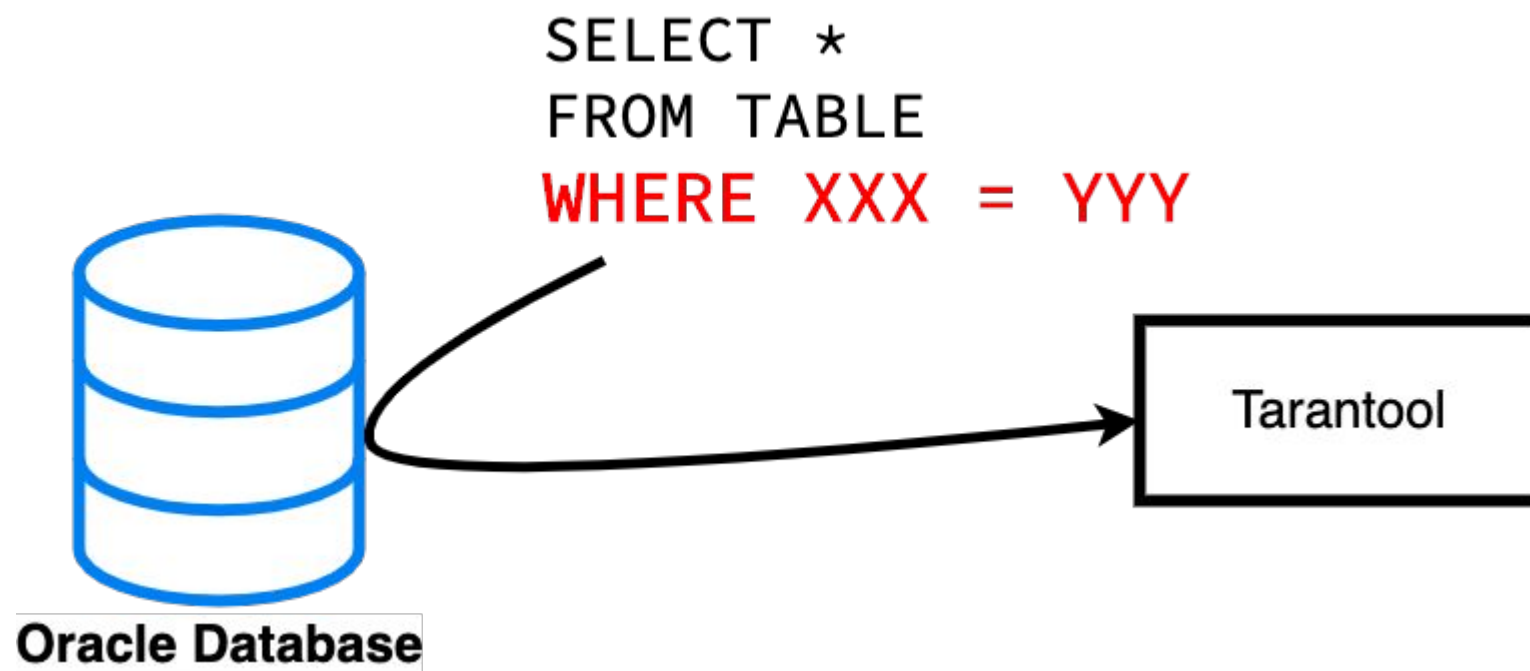


Халявы не будет

Витрина продуктов абонента:

- ~30Тb сырых данных в Oracle, ~3000 TPS обновлений
- Нет брокера, из которого можно прочесть эти изменения
- Данные в витрине должны быть транзакционно целостные
- Отставание витрины — не более 30 секунд от Oracle
- ~1Тb данных в Tarantool, до 1Тb изменений в сутки

Прогрев



Халявы не будет. Совсем

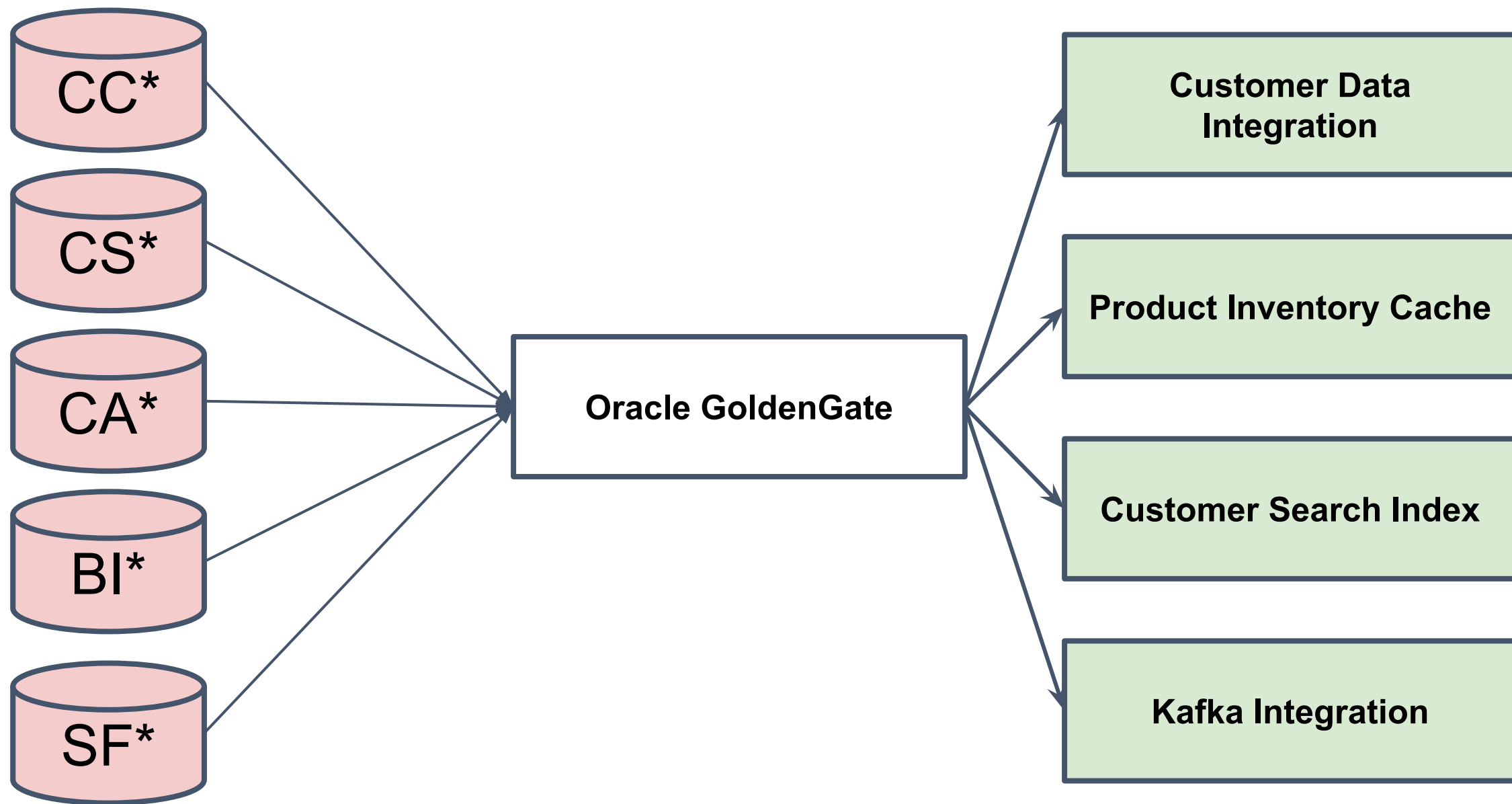
- при прогреве использовать **WHERE** нельзя
 - нагрузили Oracle
 - время выгрузки ~**2 недели**
- данные на момент выгрузки уже неактуальные
- ... и это мы еще до получения изменений не дошли

Халявы не будет. Окончательно

- Replica
- CQN
- GoldenGate

GoldenGate

- Транзакционная целостность
- Полноценная поддержка ETL
- Минимальная нагрузка на базу
- Экспертиза DBA

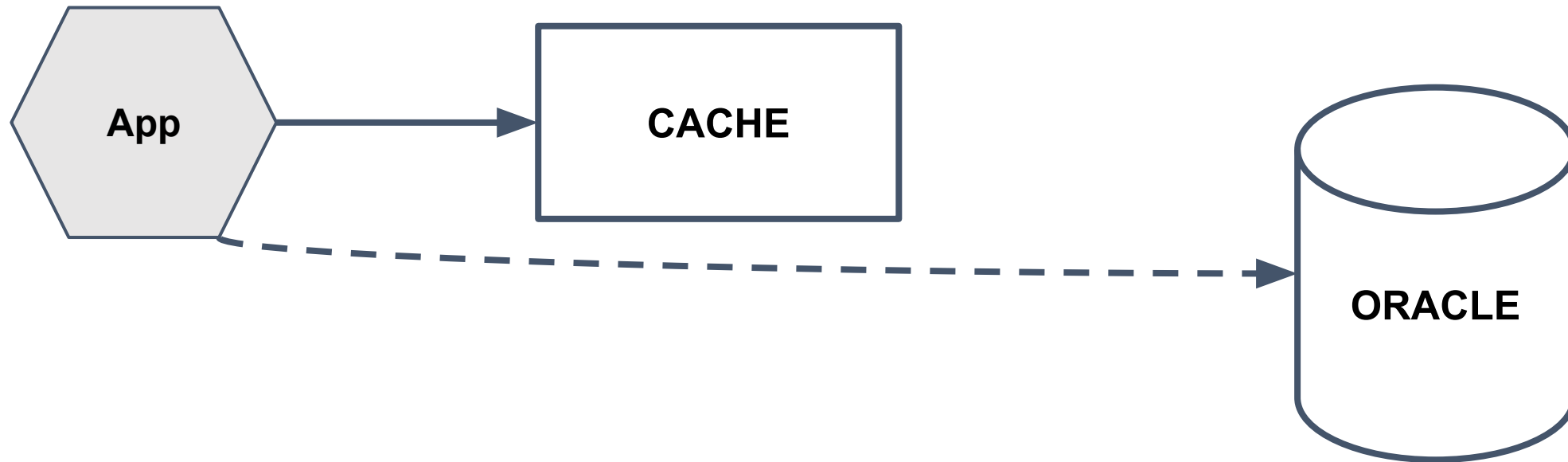


Как делали

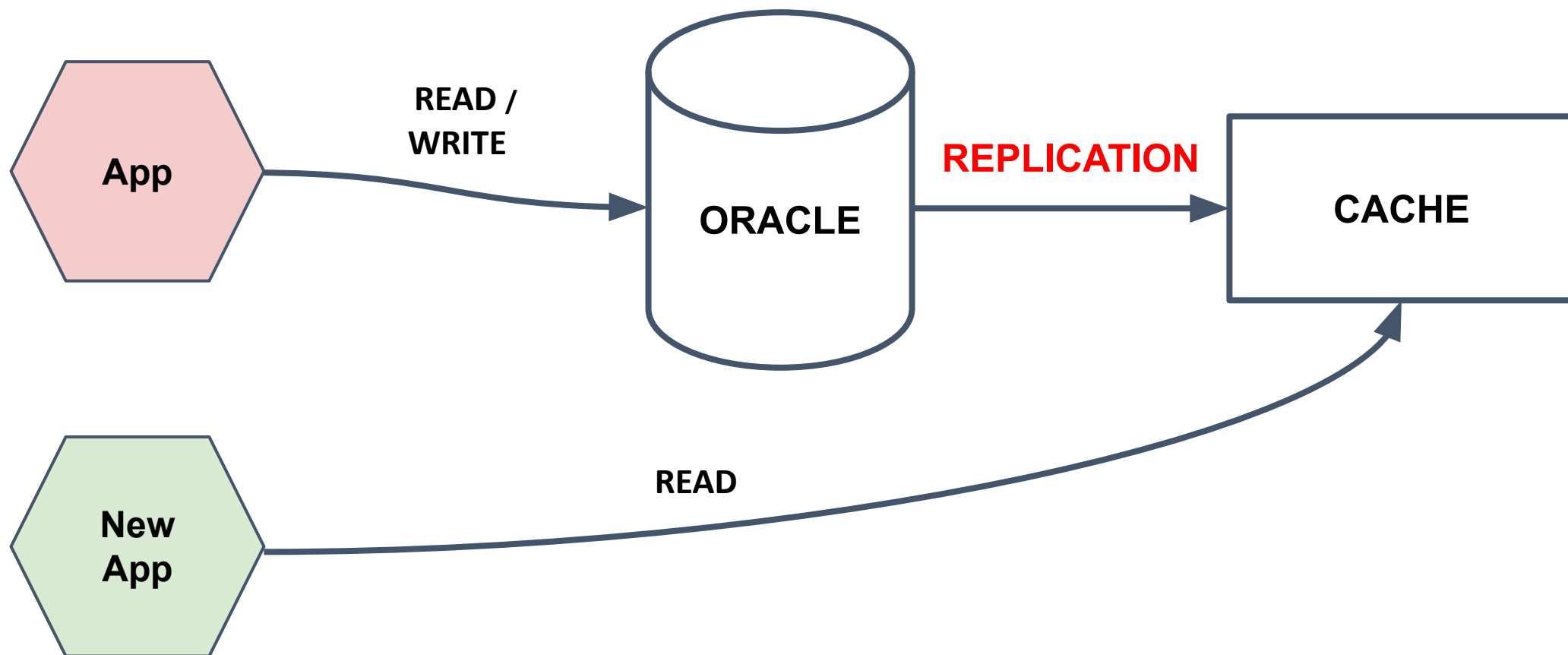
- Первый подход к снаряду:
 - Прогрев через выгрузку/загрузку CSV
 - Репликация через XML
- Второй подход к снаряду:
 - Прогрев напрямую из Oracle
 - USEREXIT
- Третий подход?....

Техника

Read through



Кэш сбоку



Работающий кэш

- **Загружен данными (прогрев)**
- **Данные актуальны (репликация)**

Работающий кэш

- Загружен данными (прогрев)
- **Данные актуальны (репликация)**

Работающий кэш

- Загружен данными (прогрев)
- Данные актуальны (репликация)
- **Observable/maintainable/*****able**

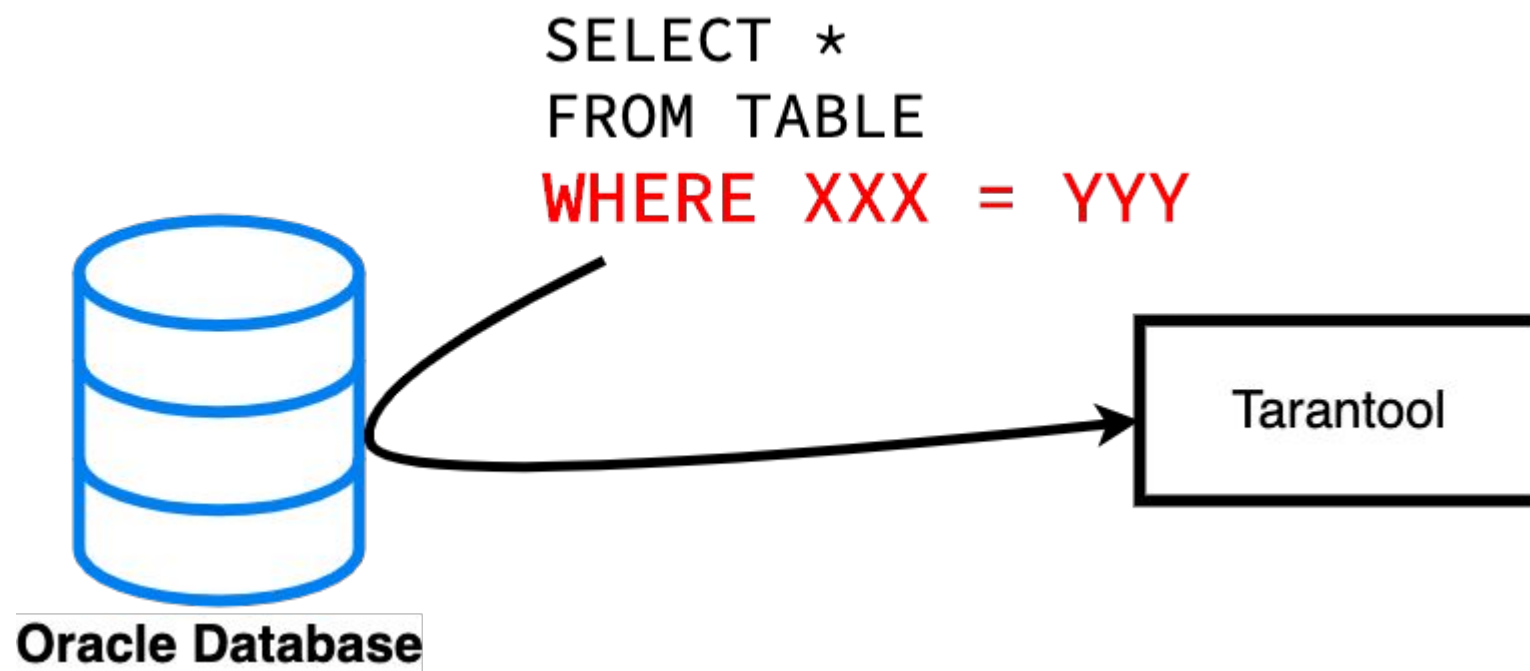
Прогрев с выгрузкой в CSV

Выгрузка из Oracle

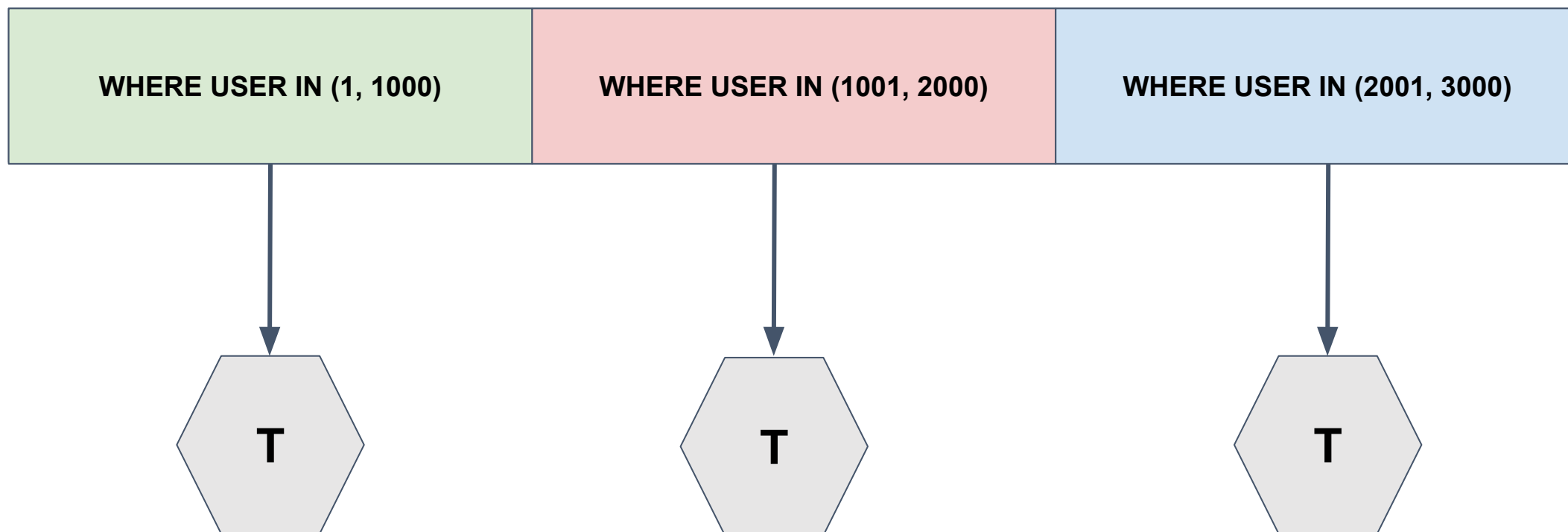
Прогрев

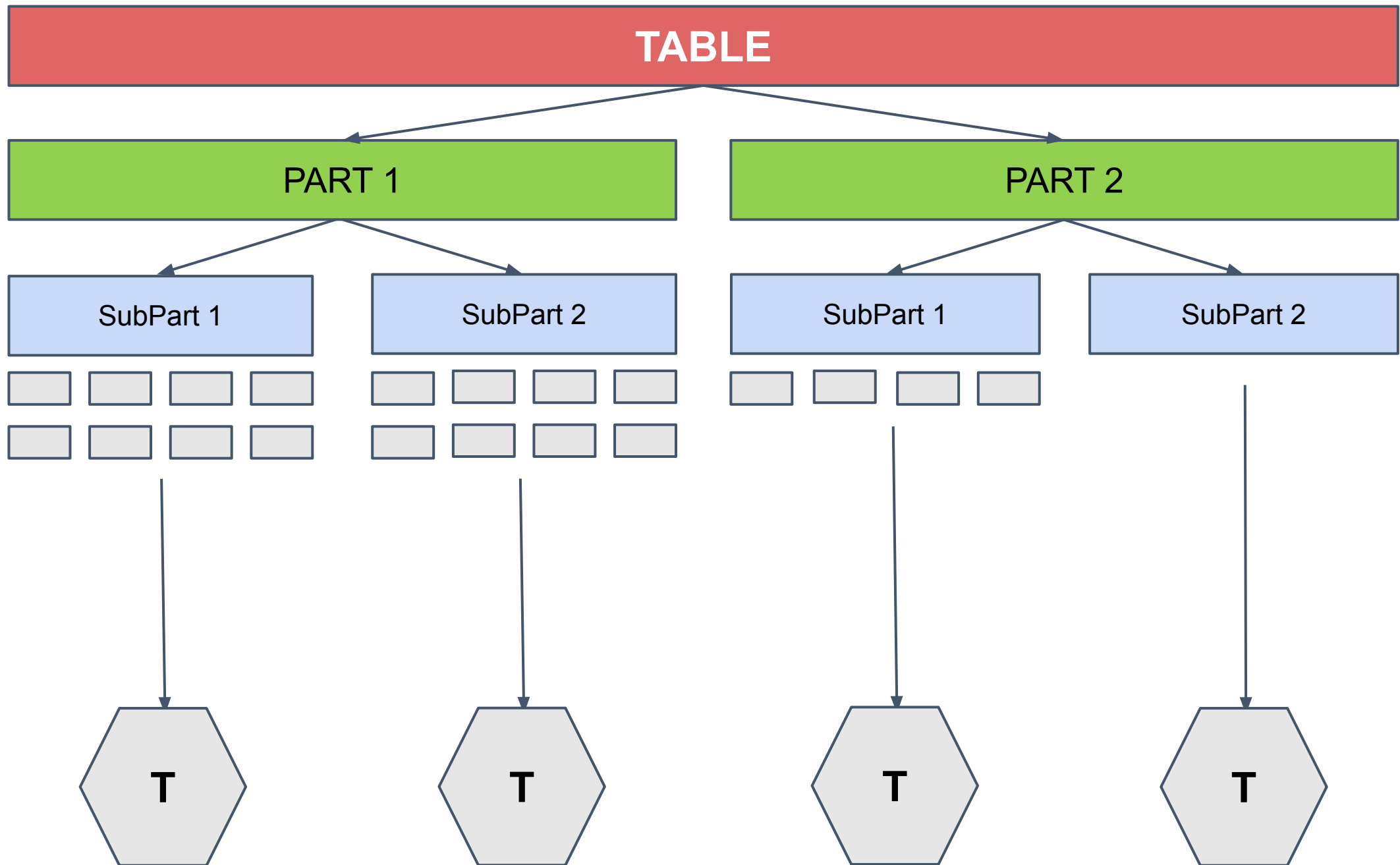
- нам нужен только срез данных
- минимальная нагрузка на Oracle
- адекватное время прогрева

Прогрев

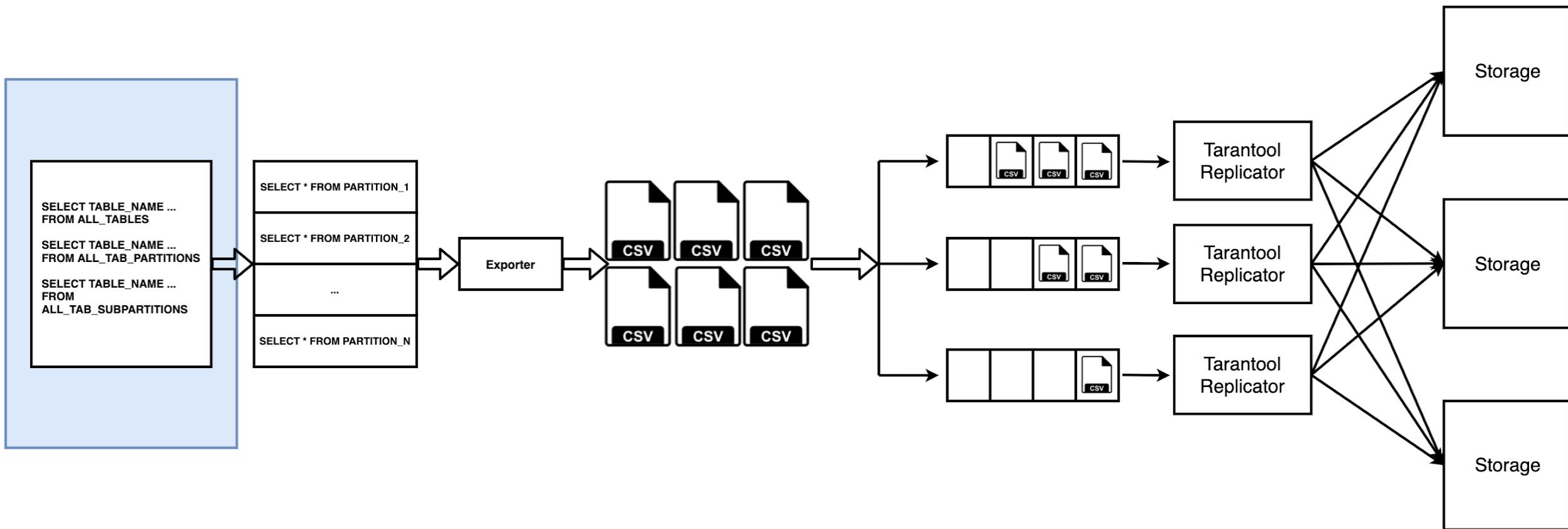


“Угнать за 60 секунд”. Минут? Часов?!

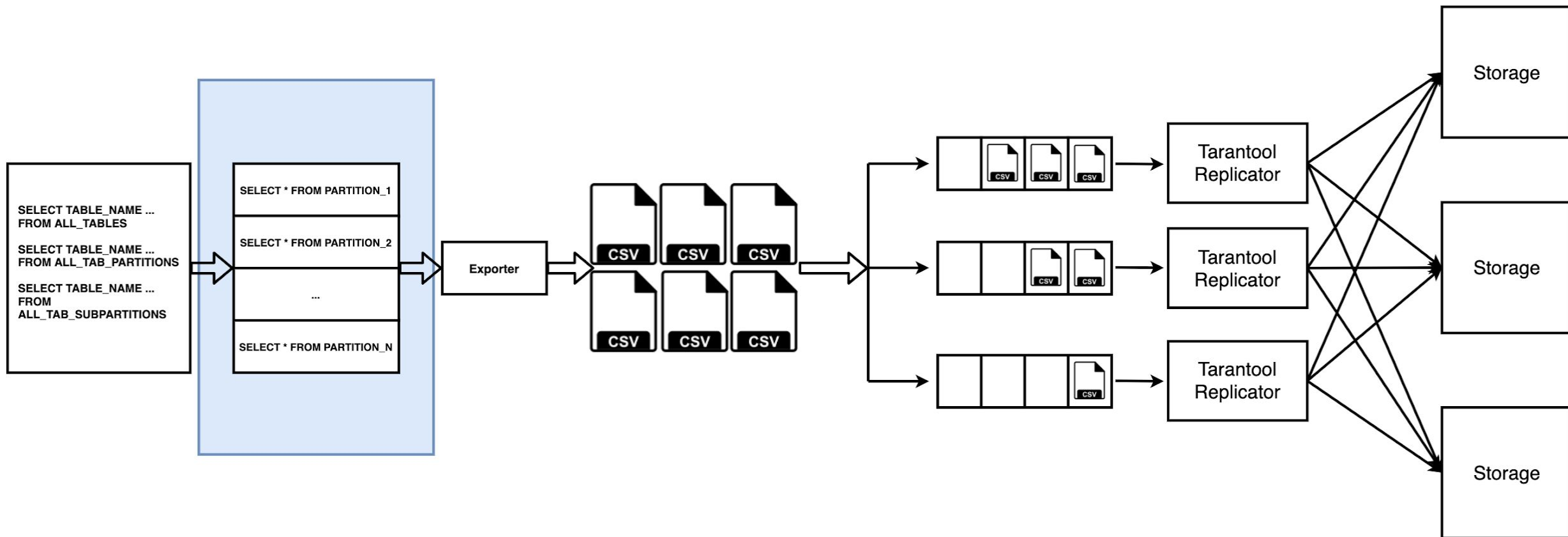




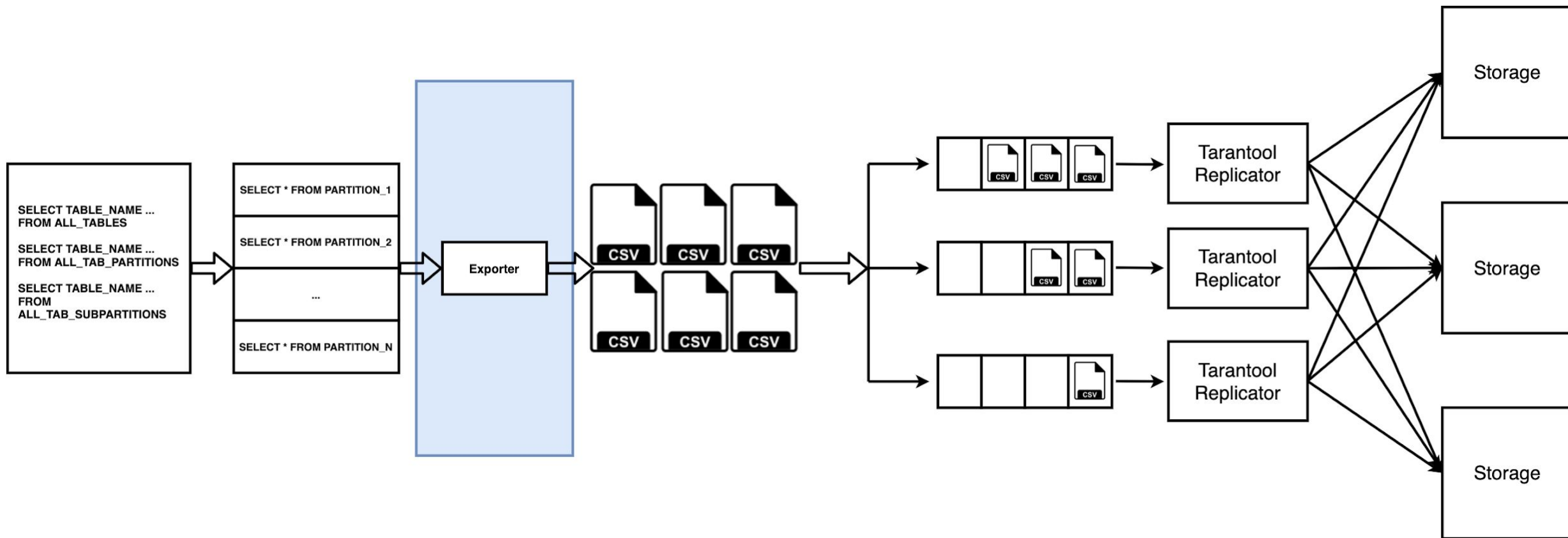
Прогрев кэша: выгрузка данных



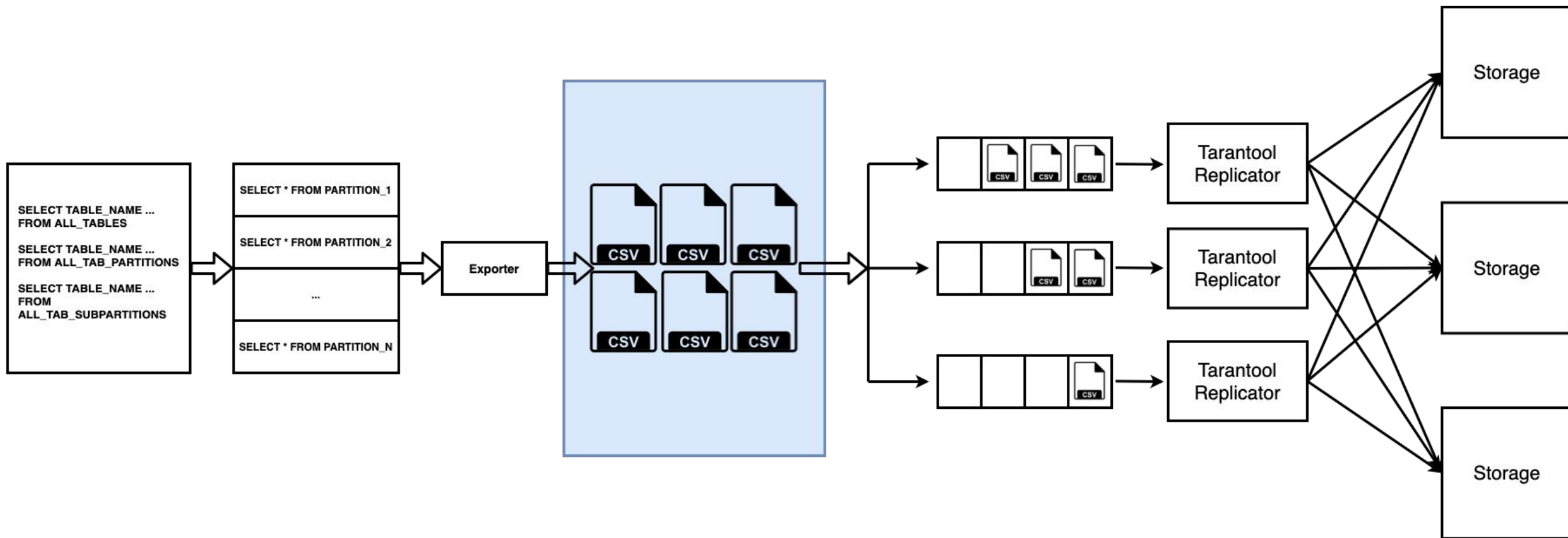
Прогрев кэша: выгрузка данных



Прогрев кэша: выгрузка данных



Прогрев кэша: выгрузка данных

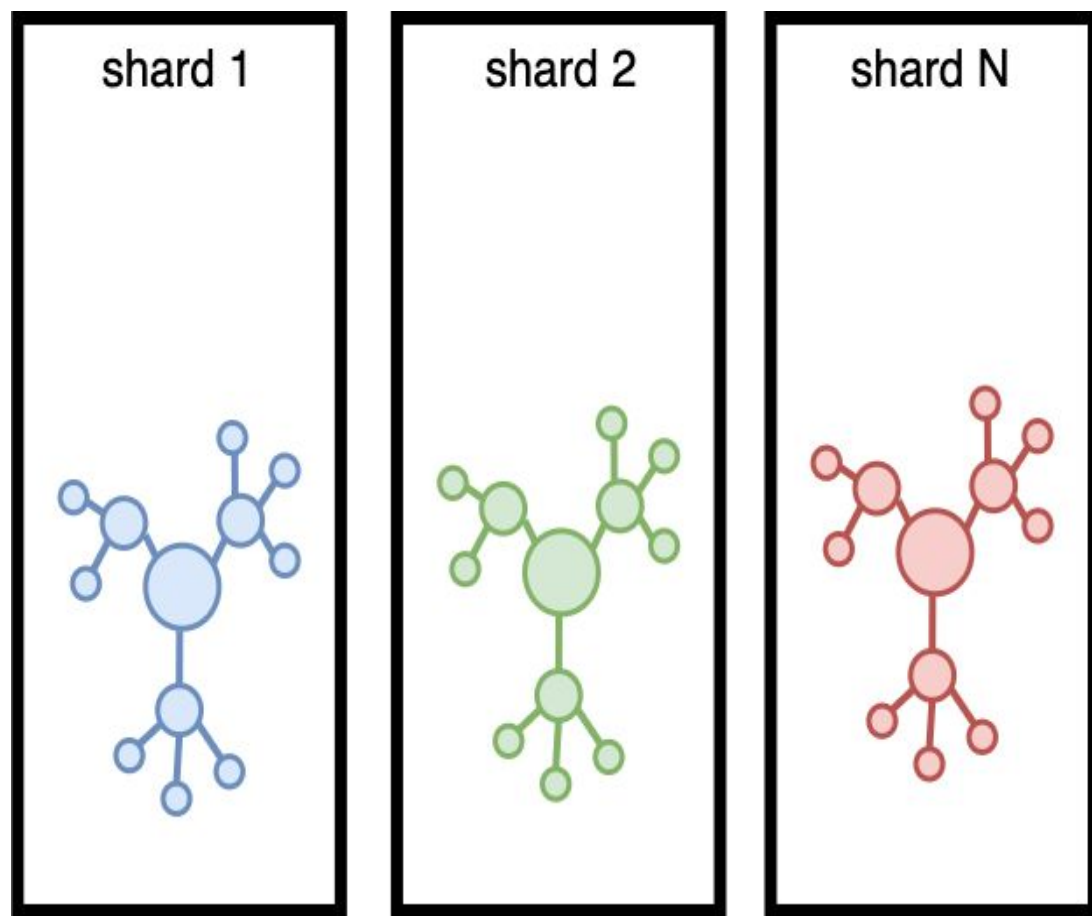


Особенности выгрузки из Oracle

- Никаких WHERE
- Адекватное партиционирование
- prefetch = 1000 (bulk)
- 10, максимум 30 потоков выгрузки
- Время выгрузки — до **6 часов**
- Скорость выгрузки — **60000 - 700 000** строк в секунду
- **1.2Tb** CSV-файлов

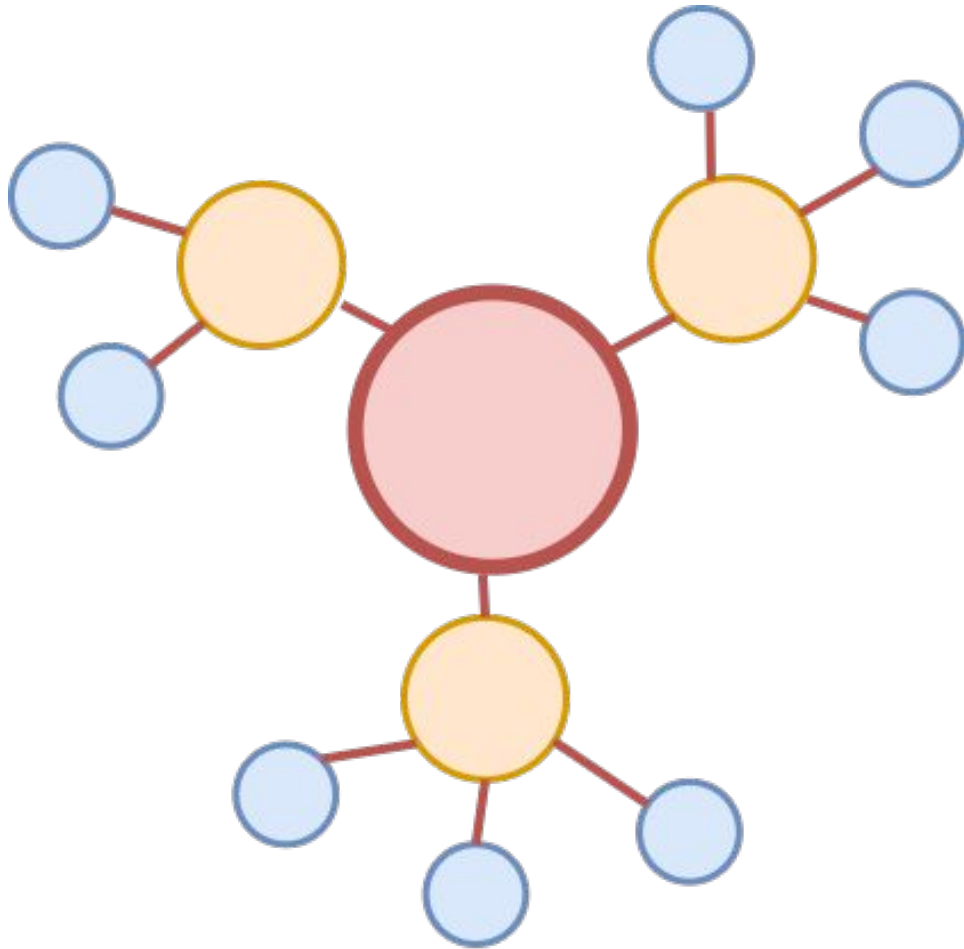
Загрузка в витрину

Модель хранения



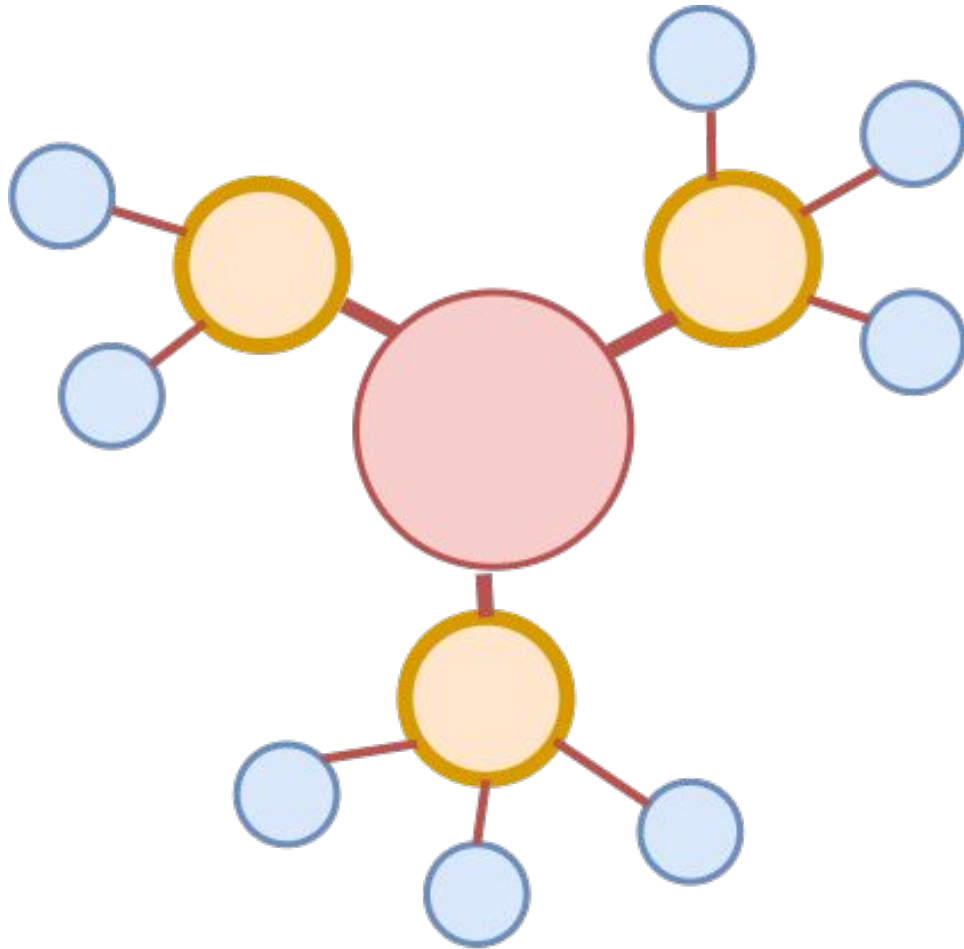
- **Данные шардированы по абоненту**
- У абонента есть подключенные услуги (~4 млрд записей)
- Для каждой услуги есть дополнительные данные (еще ~12 млрд записей)

Модель данных



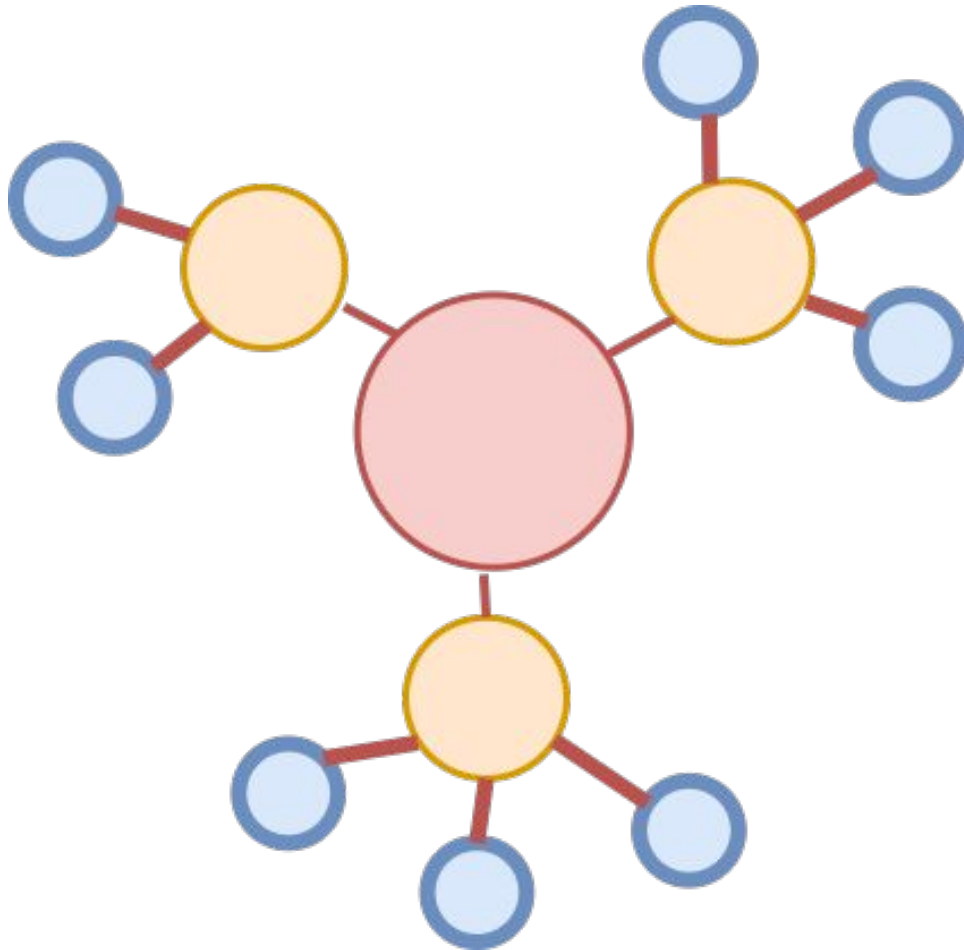
- **Данные шардированы по абоненту**
- У абонента есть подключенные услуги (~4 млрд записей)
- Для каждой услуги есть дополнительные данные (еще ~12 млрд записей)

Модель данных



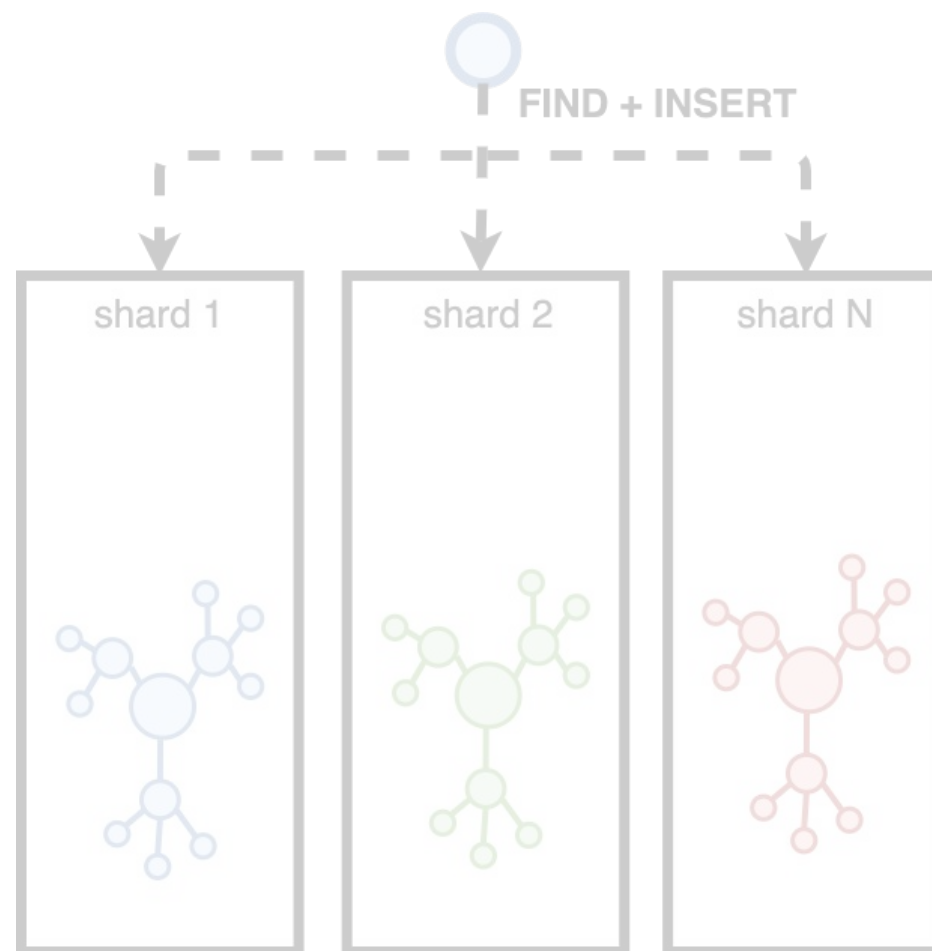
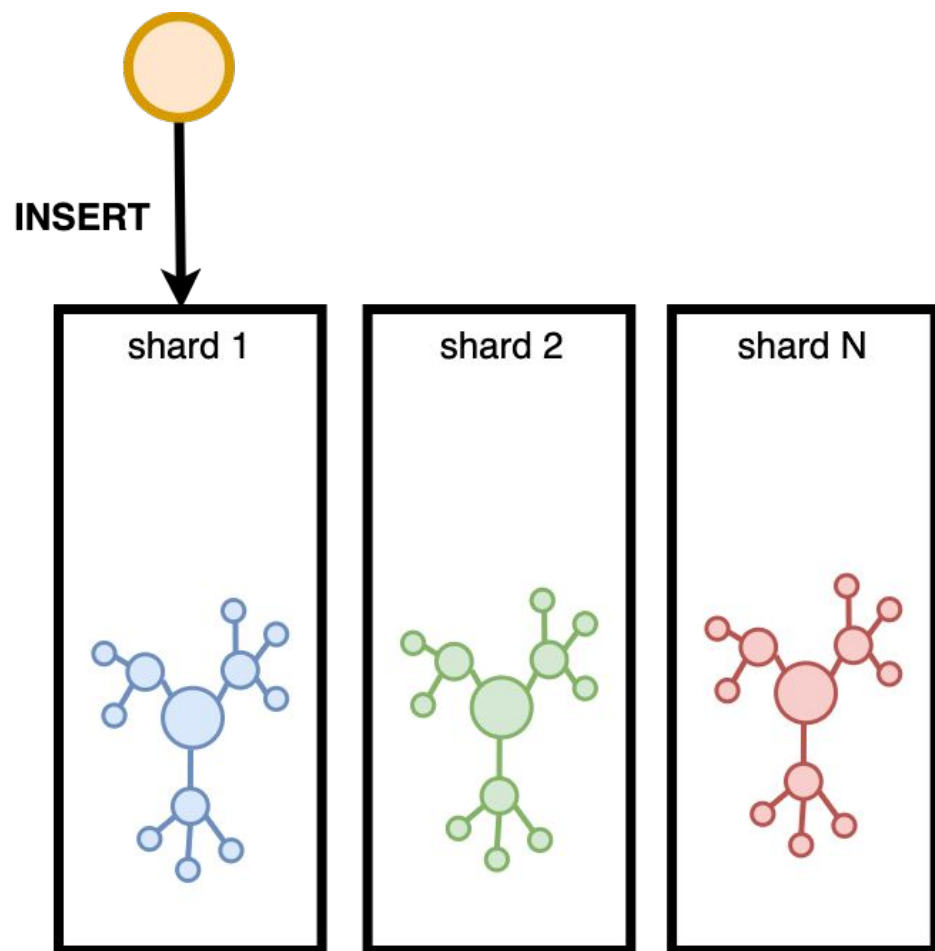
- Данные шардированы по абоненту
- **У абонента есть подключенные услуги (~4 млрд записей)**
- Для каждой услуги есть дополнительные данные (еще ~12 млрд записей)

Модель данных

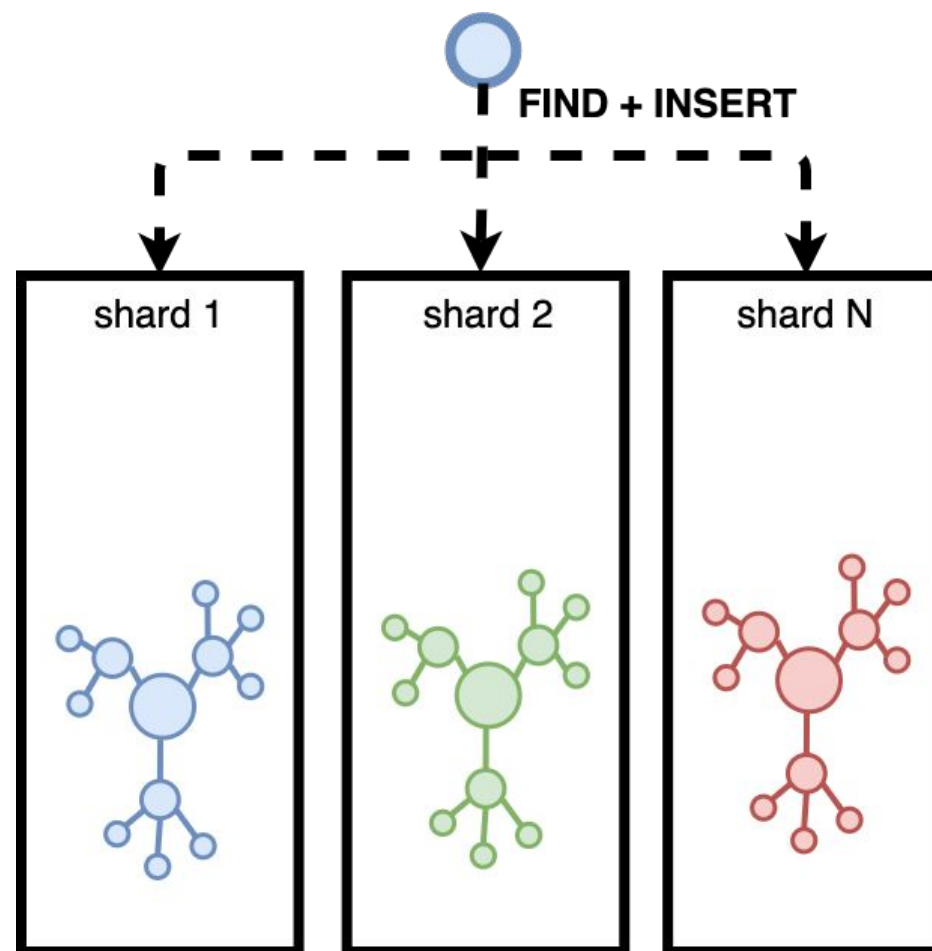
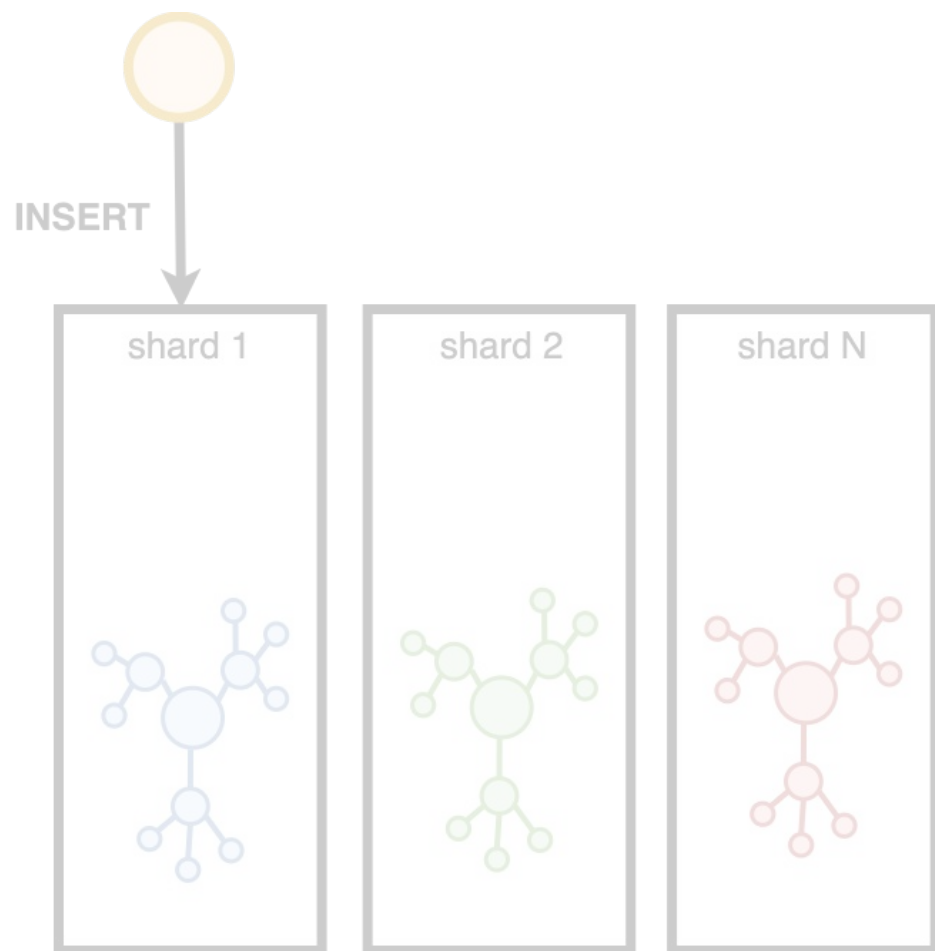


- Данные шардированы по абоненту
- У абонента есть подключенные услуги (~4 млрд записей)
- **Для каждой услуги есть дополнительные данные (еще ~12 млрд записей)**

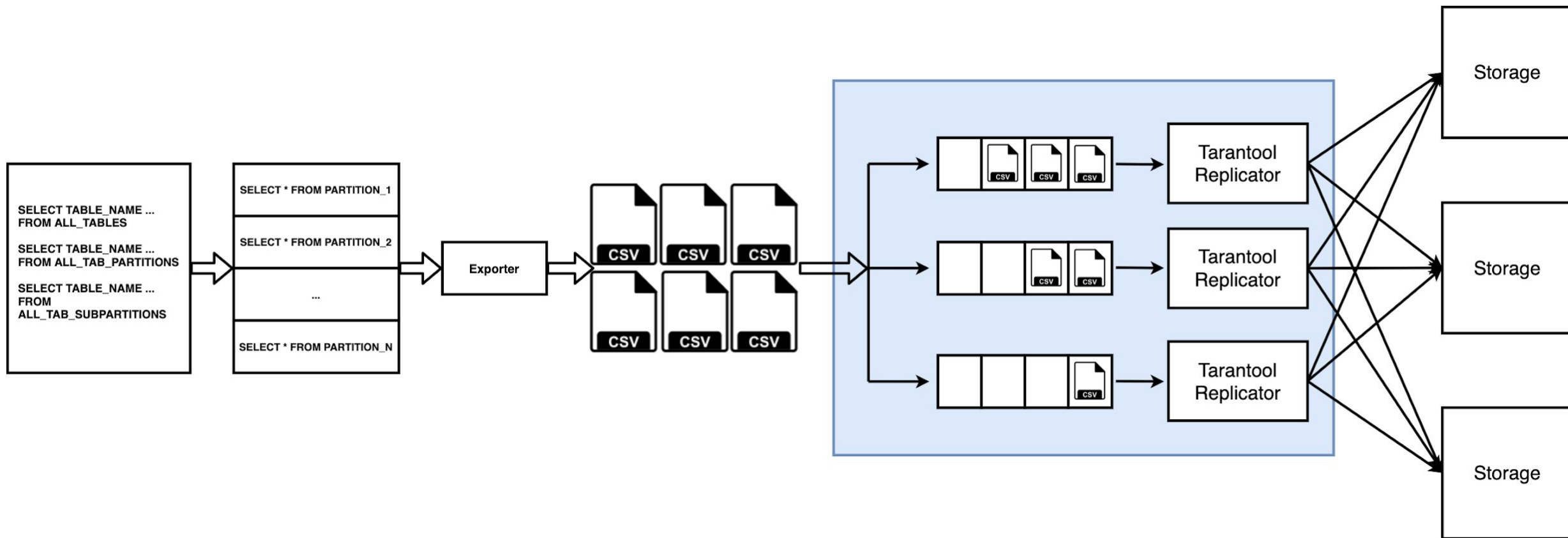
Запись данных в кэш



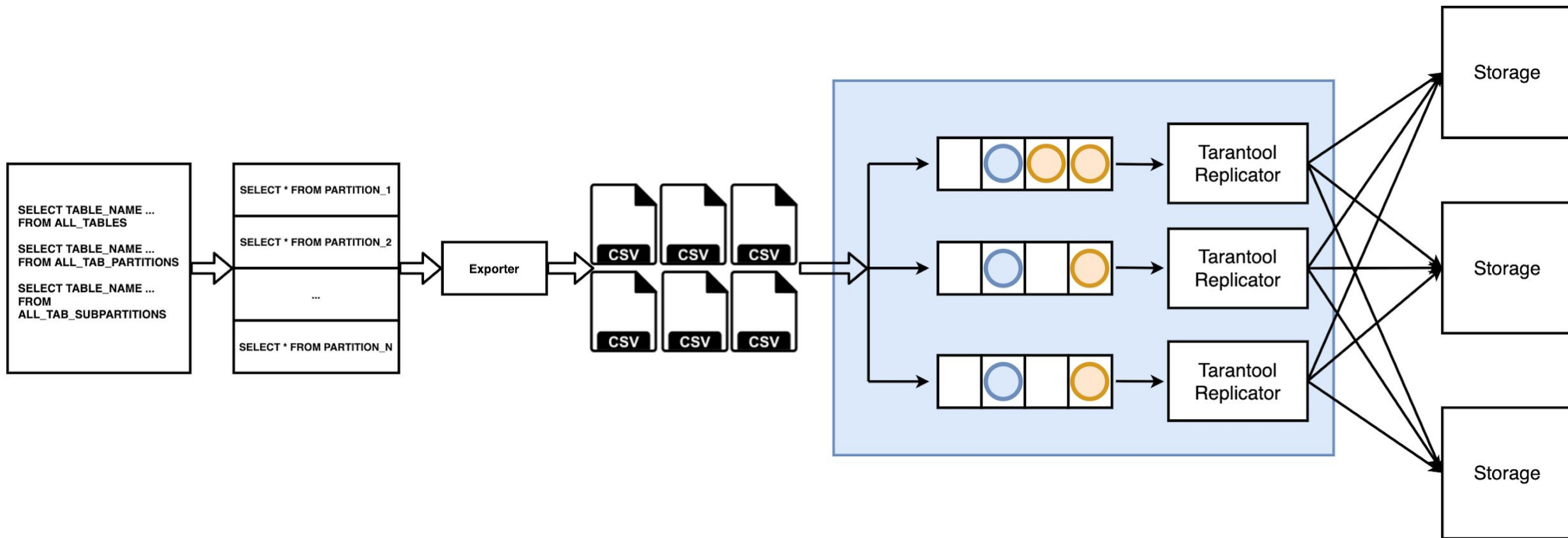
Запись данных в кэш



Прогрев кэша: загрузка данных



Прогрев кэша: порядок загрузки



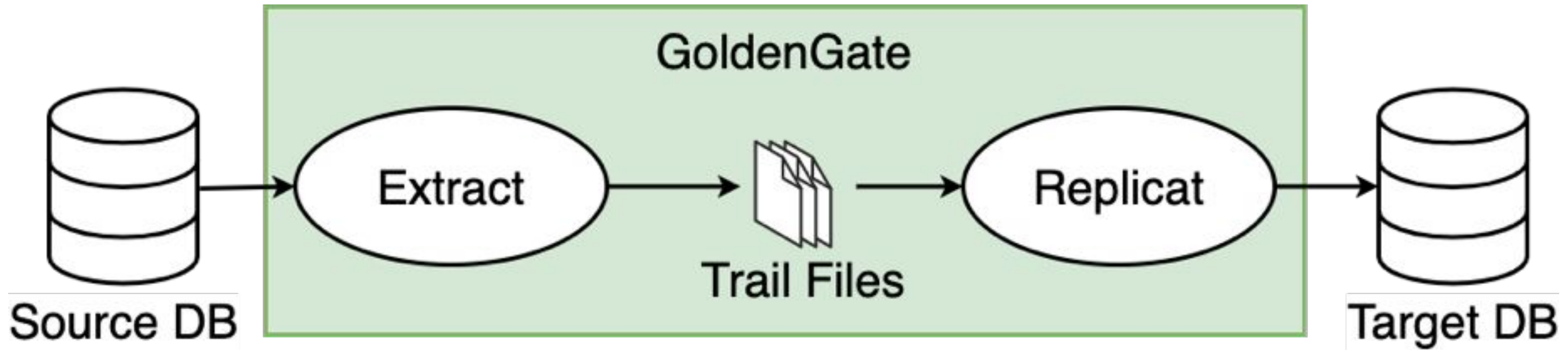
Особенности загрузки

- 1 шаг — грузим данные по абоненту
- 2 шаг — догружаем все привязанные к нему сущности
- отфильтровали данные на загрузке в кэш
- ~12 часов на загрузку данных

Change Data Capture

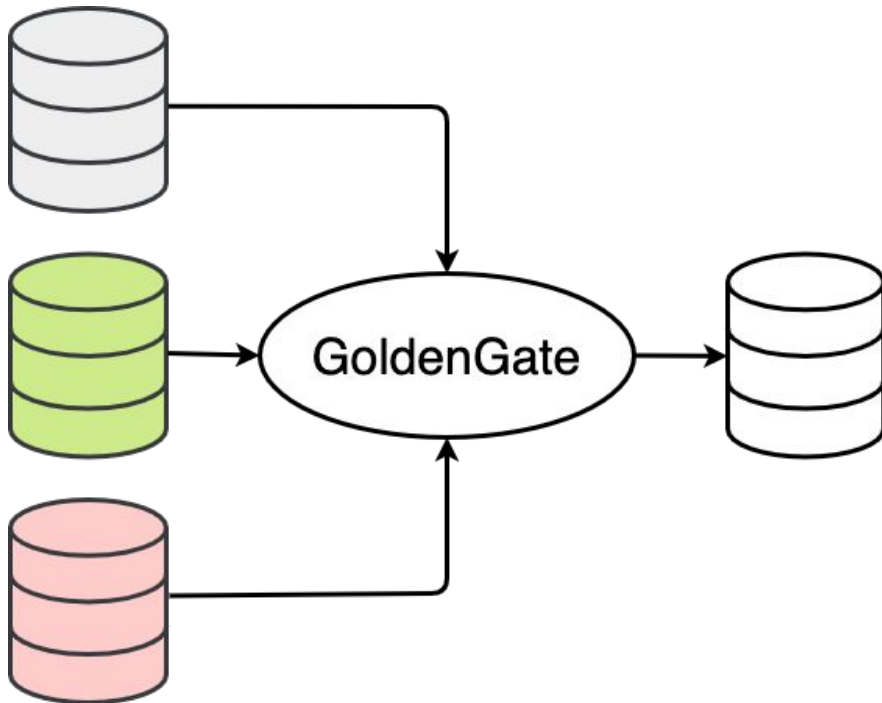
Oracle GoldenGate

GoldenGate

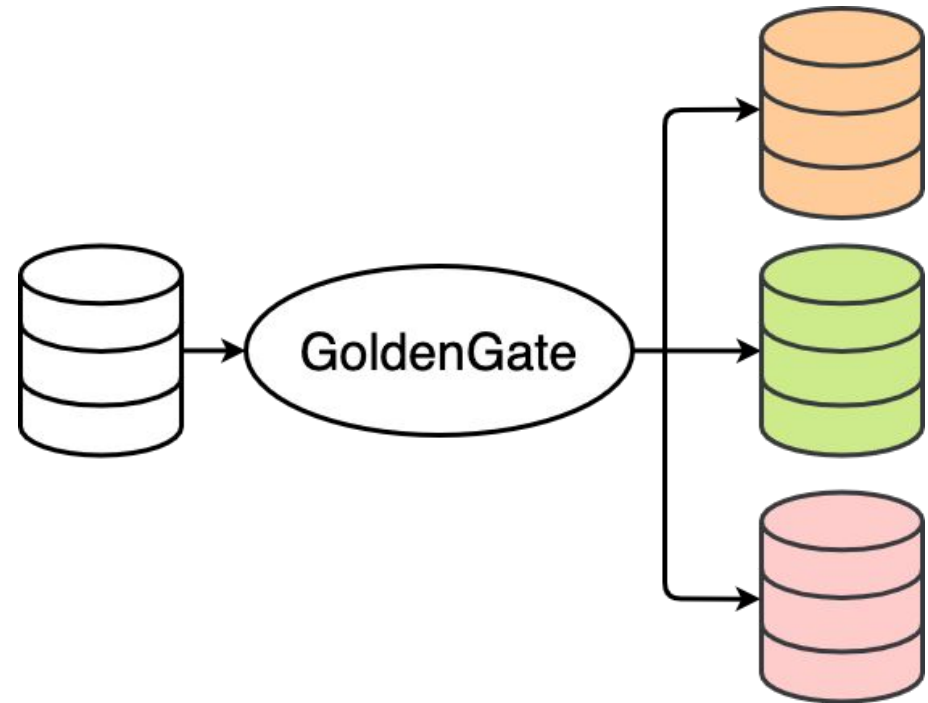


GoldenGate: различные топологии потоков данных

Consolidation

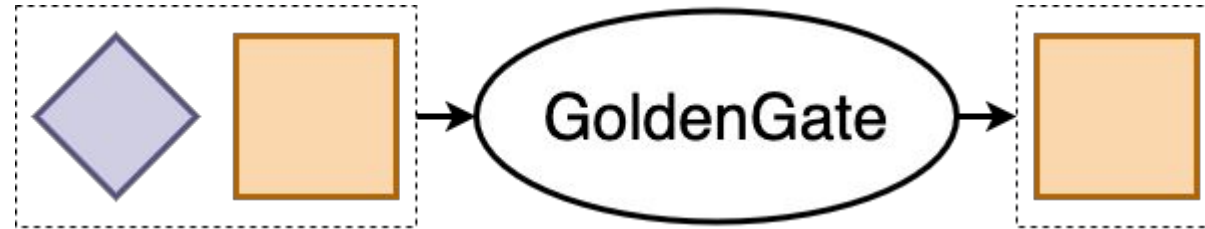


Duplication



GoldenGate: фильтрация и трансформация

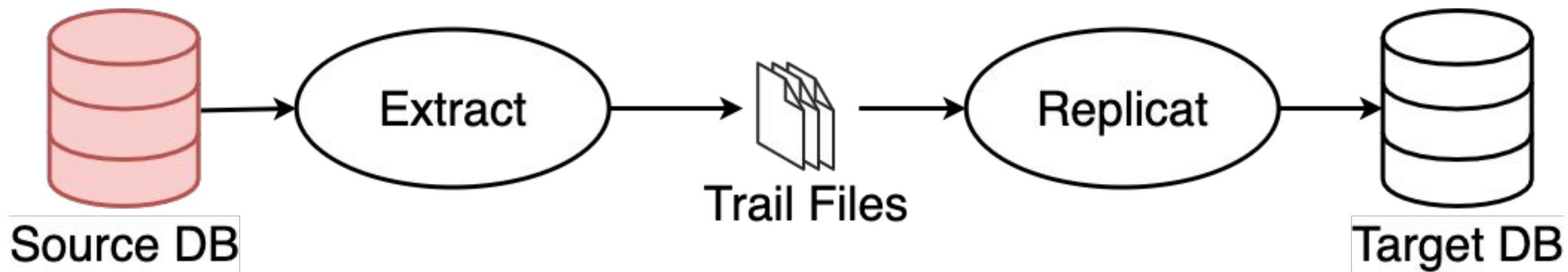
Filter:



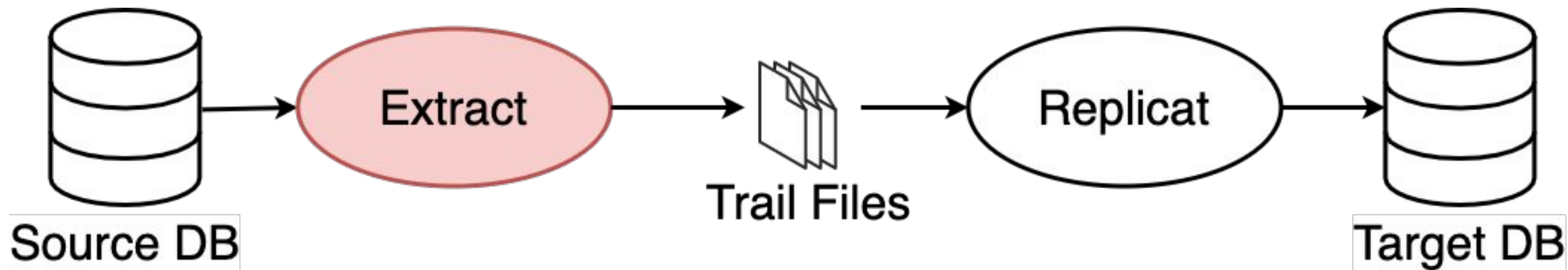
Transform:



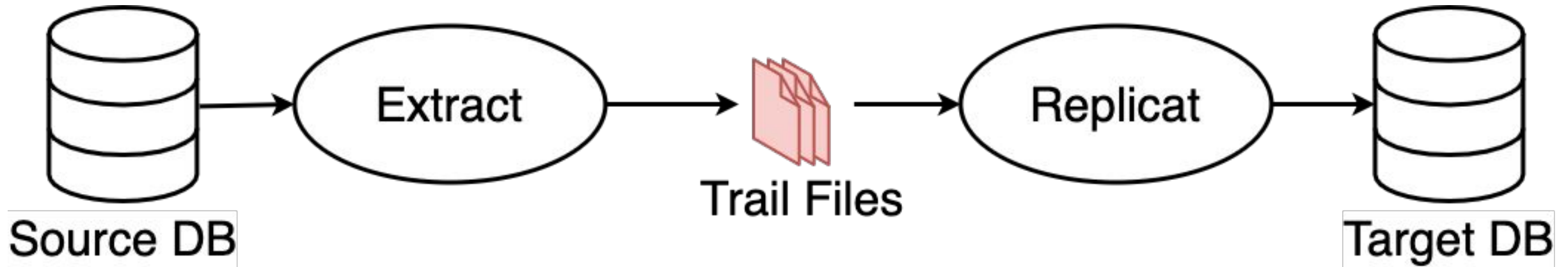
GoldenGate: как работает



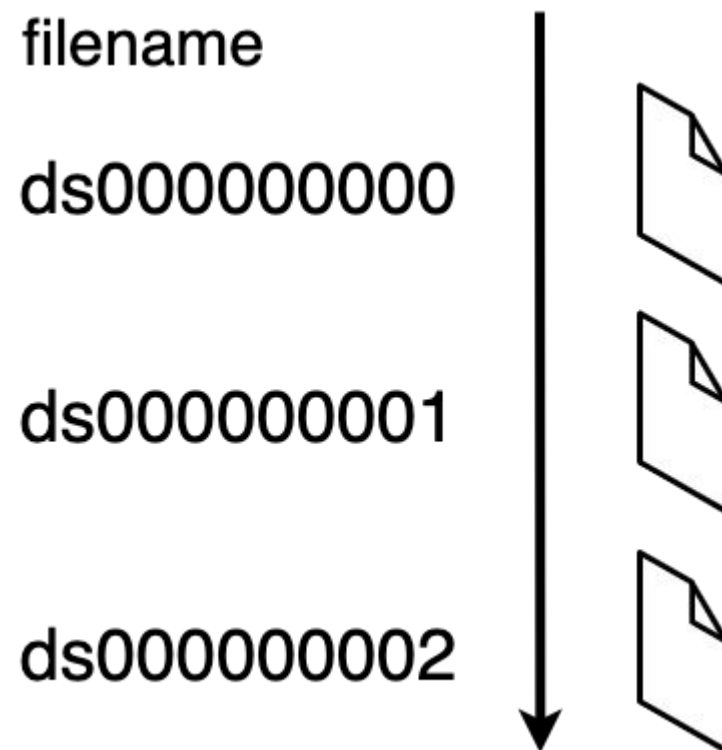
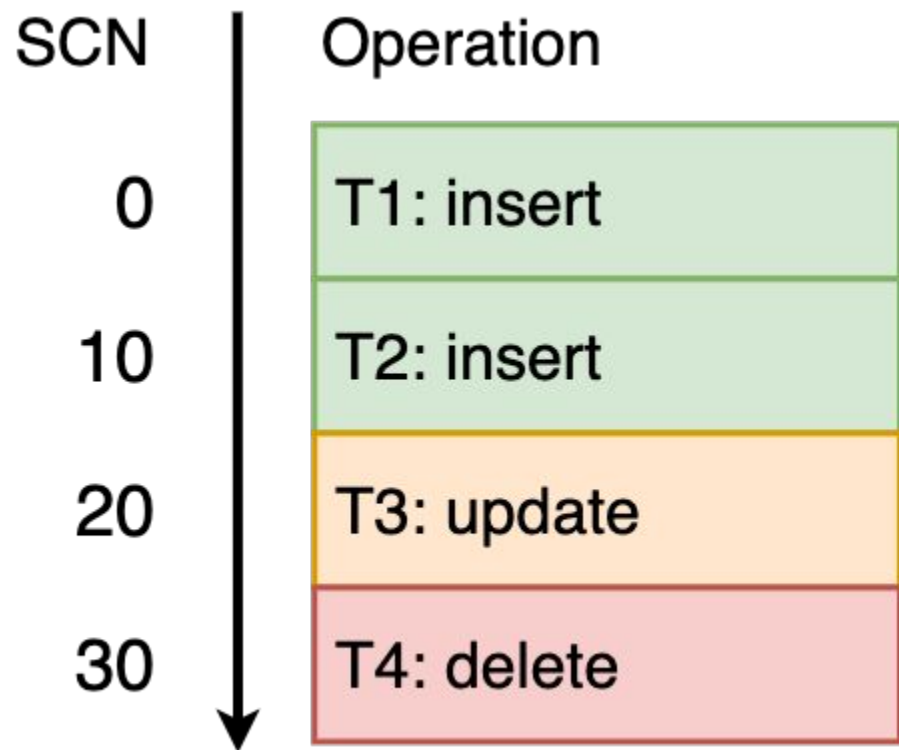
GoldenGate: как работает



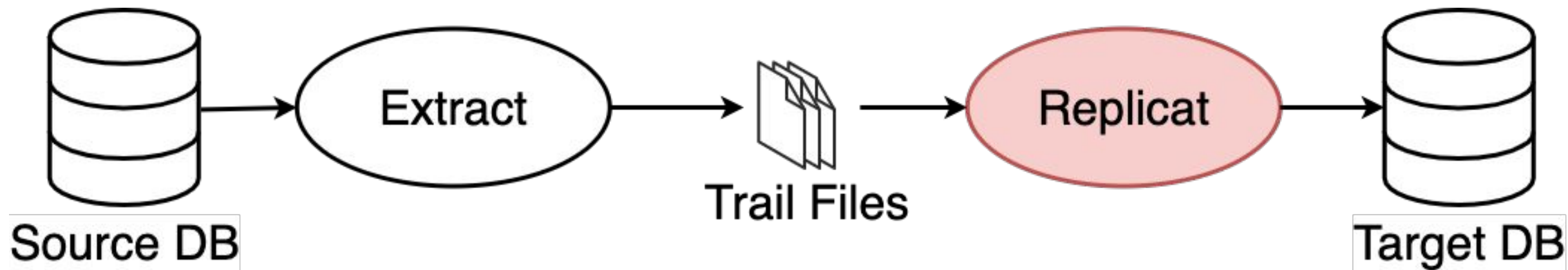
GoldenGate: как работает



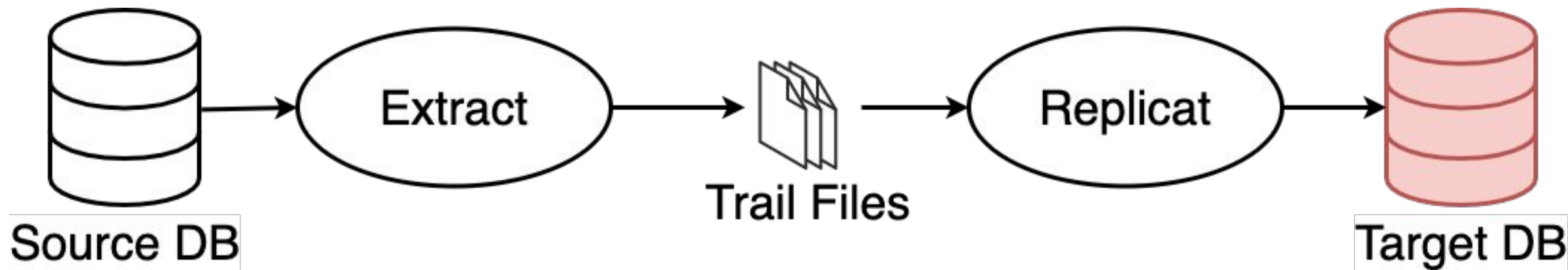
Trail: log транзакций



GoldenGate: как работает



GoldenGate: как работает



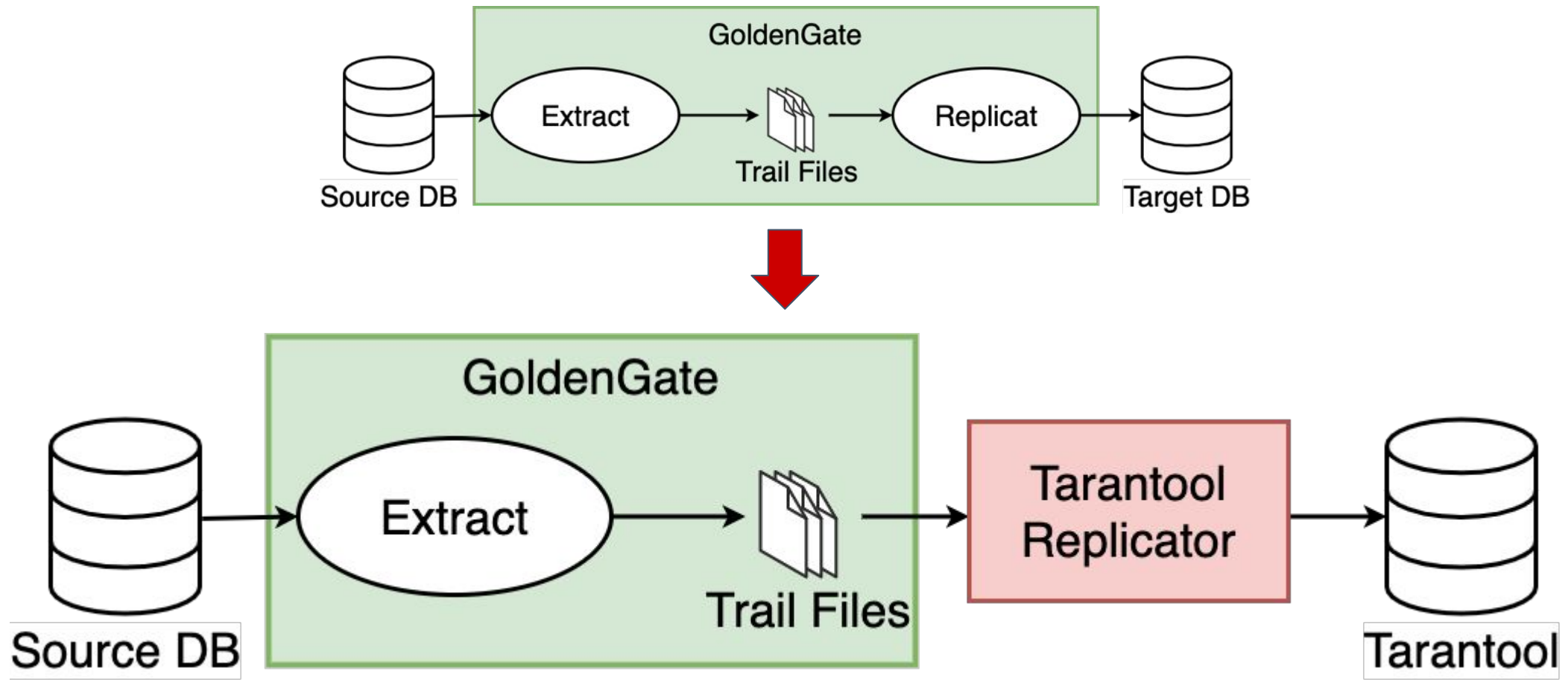
Summary

- Позволяет реплицировать данные
- Может изменять реплицируемые данные
- Реплицирует только закоммиченные транзакции

Oracle -> Tarantool

XML

GoldenGate: Tarantool



Trail: форматы

- Бинарный
- ASCII
- SQL
- **XML**

Trail: XML

```
<dbupdate table="XXX" type="update" image="before">
  <columns>
    <column name="COLUMN1">value 1</column>
  </columns>
  <tokens>
    <token name="X_COMMIT_TIME">1620802191</token>
  </tokens>
</dbupdate>
```

Trail: XML

```
<dbupdate table="XXX" type="update" image="before">
  <columns>
    <column name="COLUMN1">value 1</column>
  </columns>
  <tokens>
    <token name="X_COMMIT_TIME">1620802191</token>
  </tokens>
</dbupdate>
```

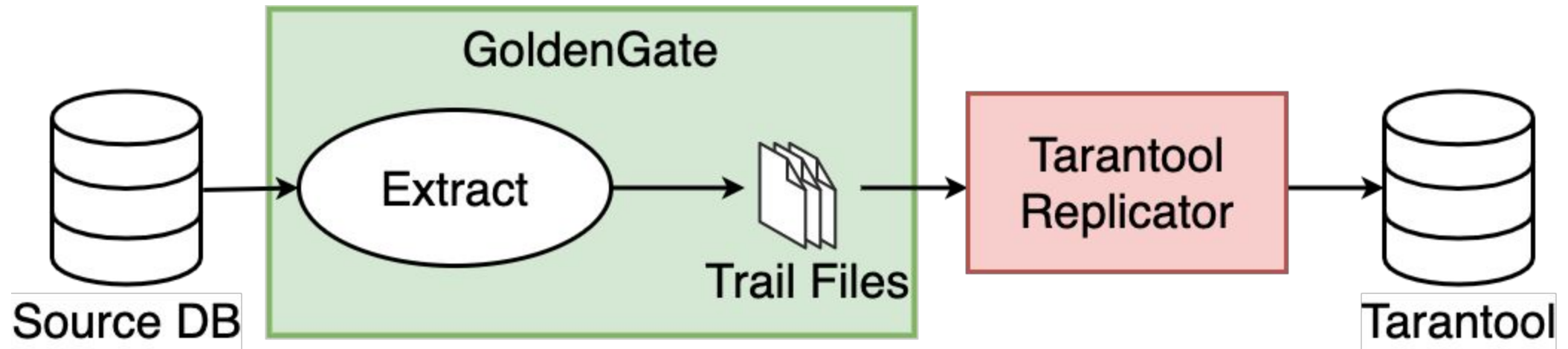
Trail: XML

```
<dbupdate table="XXX" type="update" image="before">  
  <columns>  
    <column name="COLUMN1">value 1</column>  
  </columns>  
  <tokens>  
    <token name="X_COMMIT_TIME">1620802191</token>  
  </tokens>  
</dbupdate>
```


Trail: XML

```
<dbupdate table="XXX" type="update" image="before">  
  <columns>  
    <column name="COLUMN1">value 1</column>  
  </columns>  
  <tokens>  
    <token name="X_COMMIT_TIME">1620802191</token>  
  </tokens>  
</dbupdate>
```

GoldenGate: Tarantool



Итого

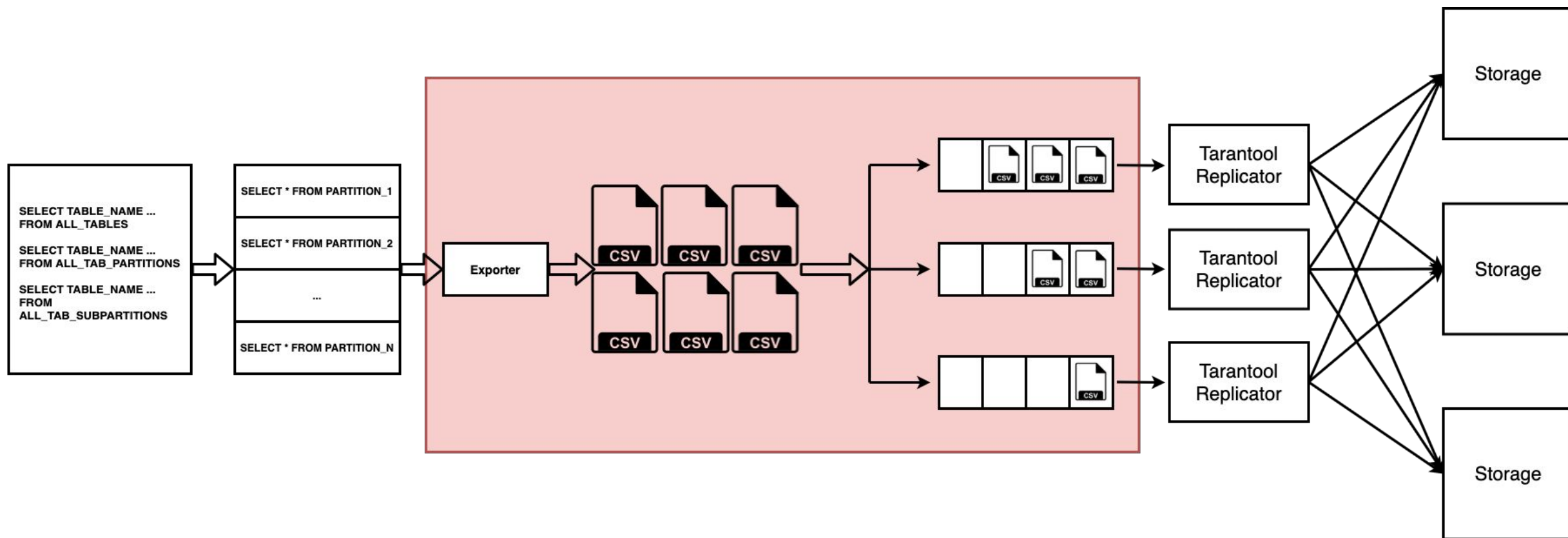
- Грузим данные из CSV
- Репликация из XML
- CSV хорошо на отладке
- XML простой: открыл файл — видишь транзакции

Round 2

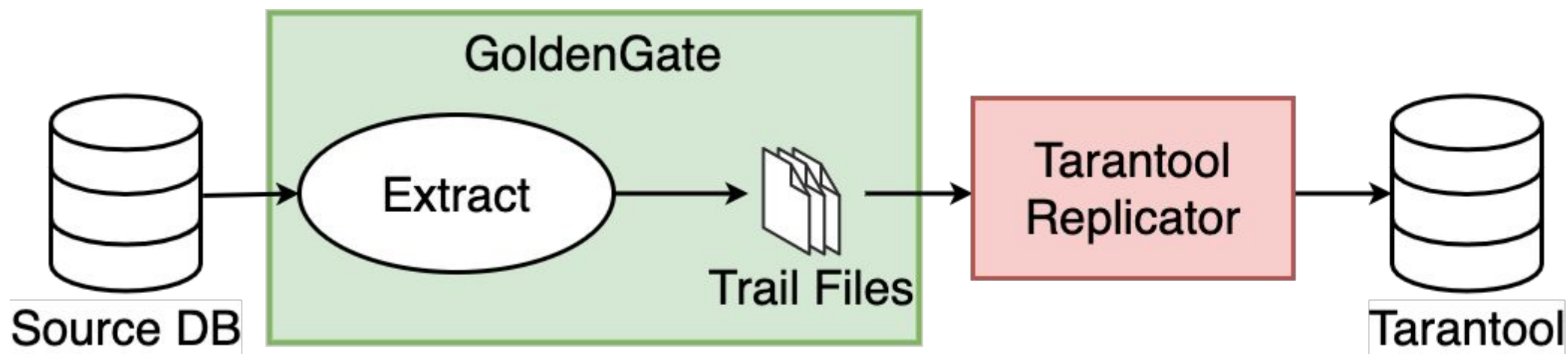
FIGHT!

Что можно сделать лучше?

Оптимизируем прогрев



Как не написать свой GG?



XML: не все типы данных

```
<dbupdate table="XXX" type="update" image="before">
```

Unsupported data type CLOB for column COLUMN2

```
  <columns>
```

```
    <column name="COLUMN1">value 1</column>
```

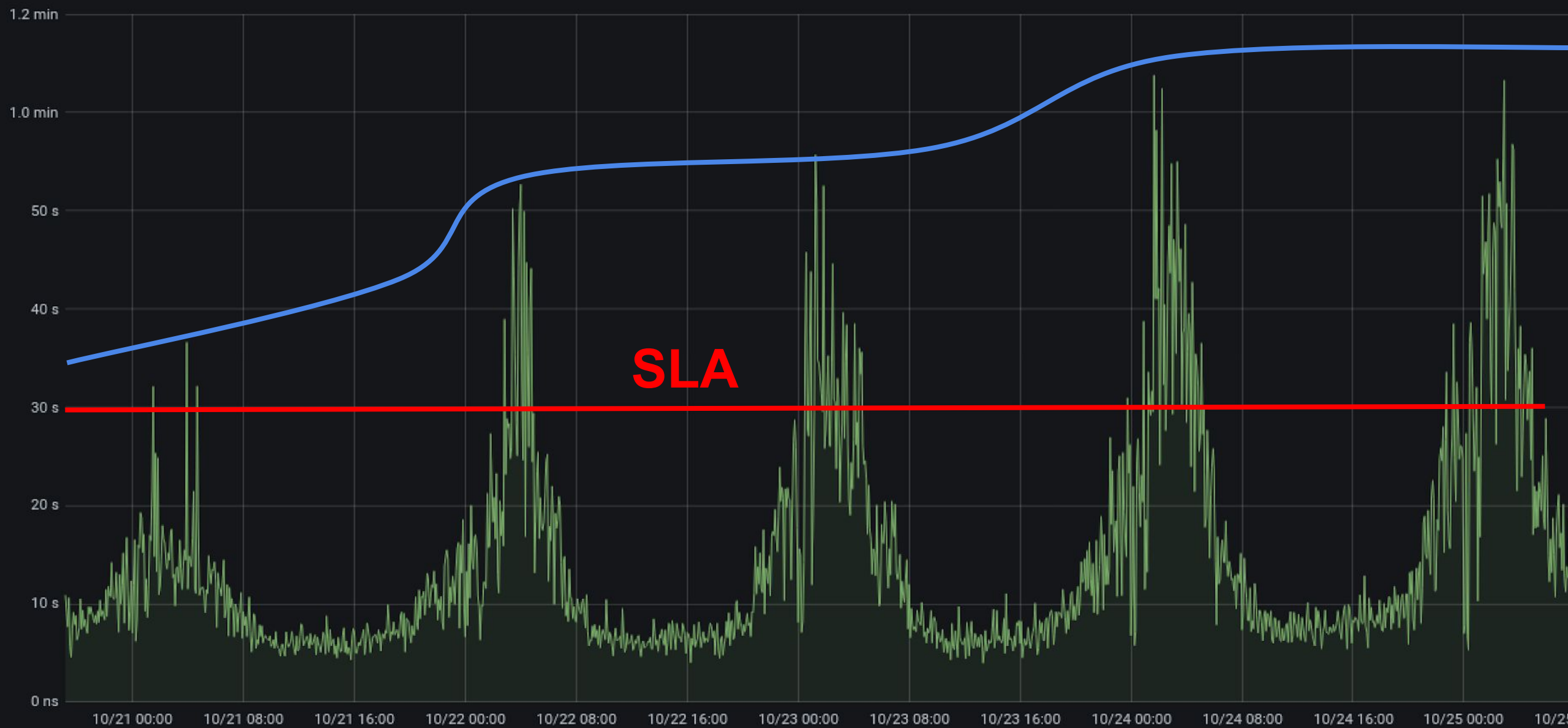
```
  </columns>
```

```
</dbupdate>
```


Оставание репликации по времени (секунды)



Оставание репликации по времени (секунды)

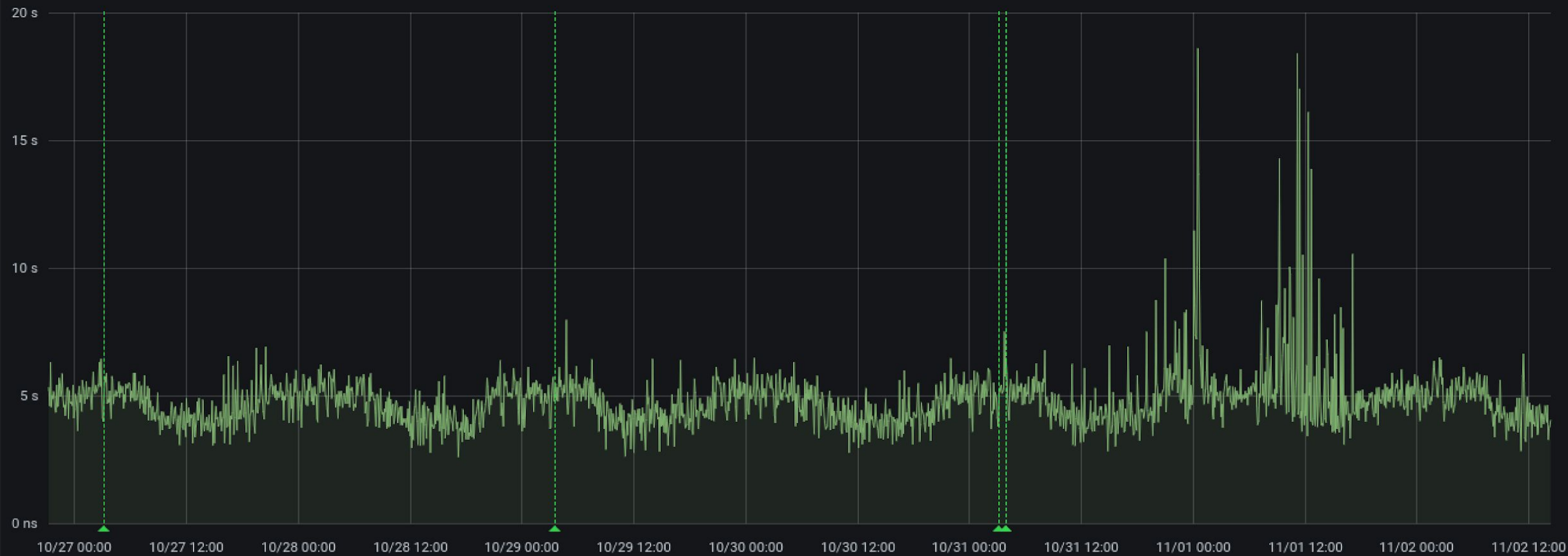


SLA

Меньше транзакций -> больше лаг

- транзакции набираем пачками
- не обрабатываем, пока не набрали пачку
- добавили **flush timeout**

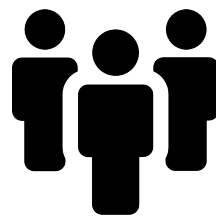
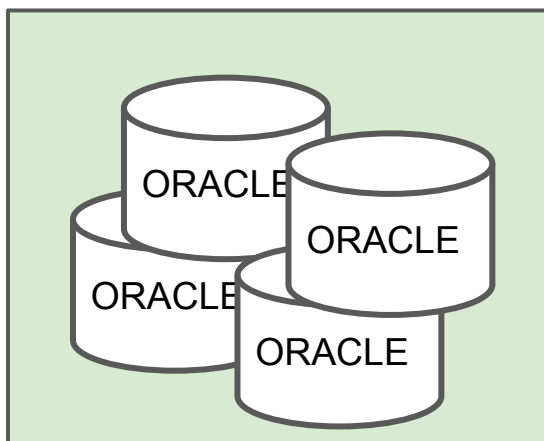
Оставание репликации по времени (секунды)



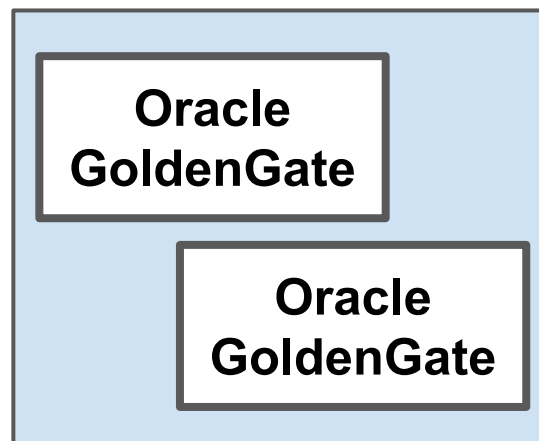
Как не ломать границы эксплуатации?



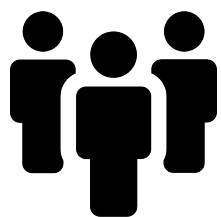
DBA



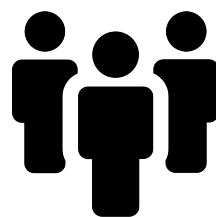
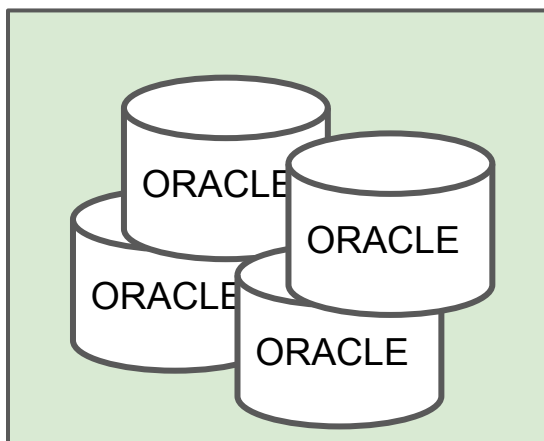
GG



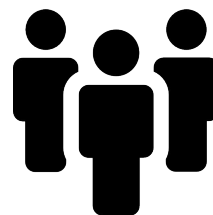
Как не ломать границы эксплуатации?



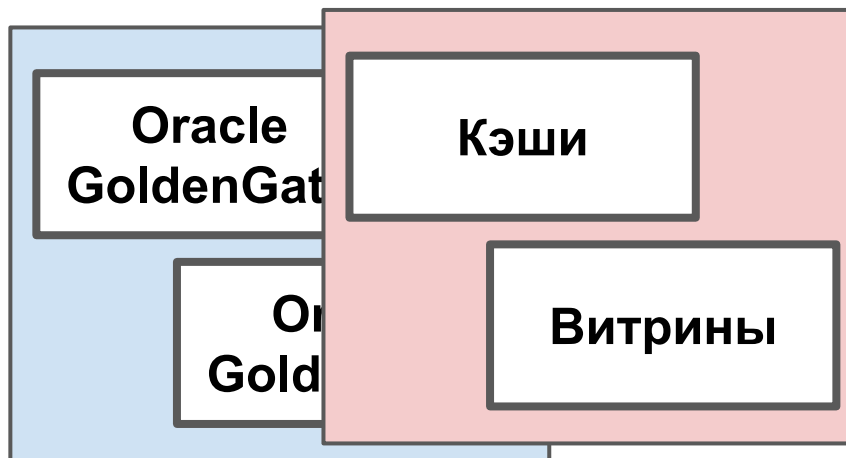
DBA



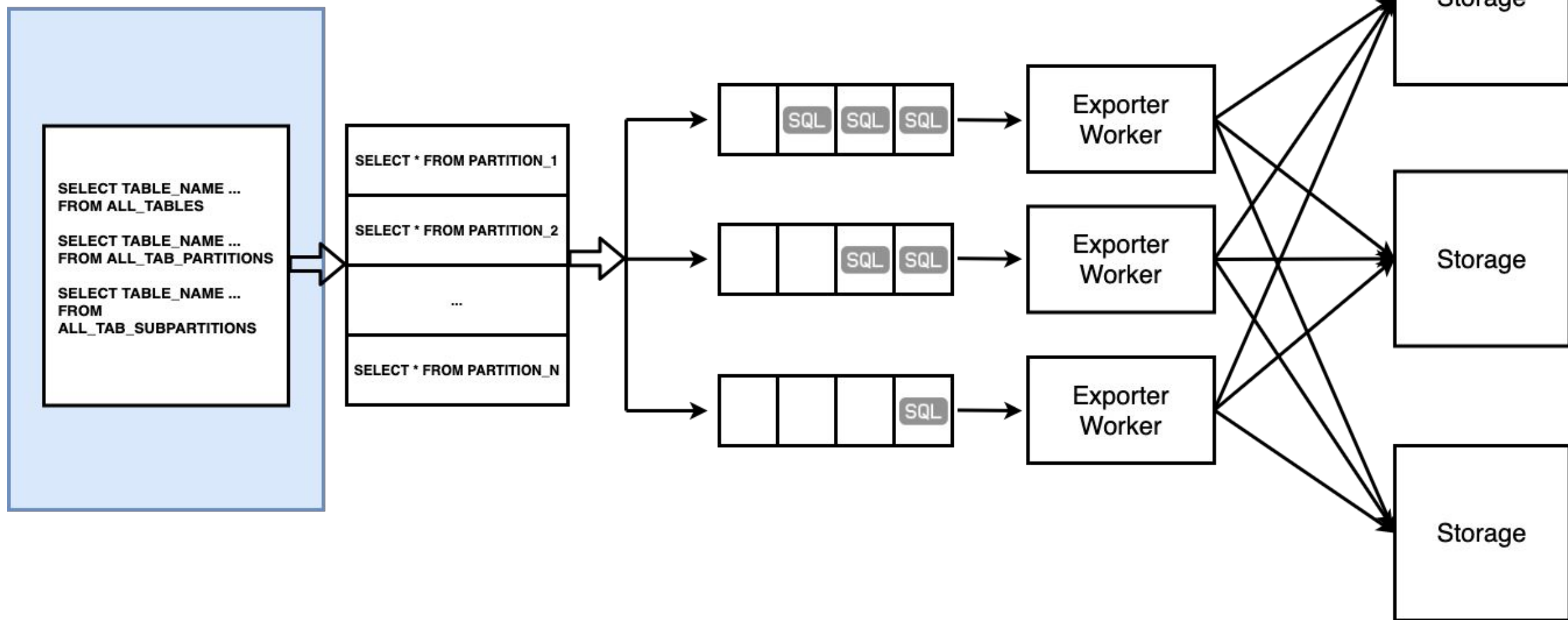
GG

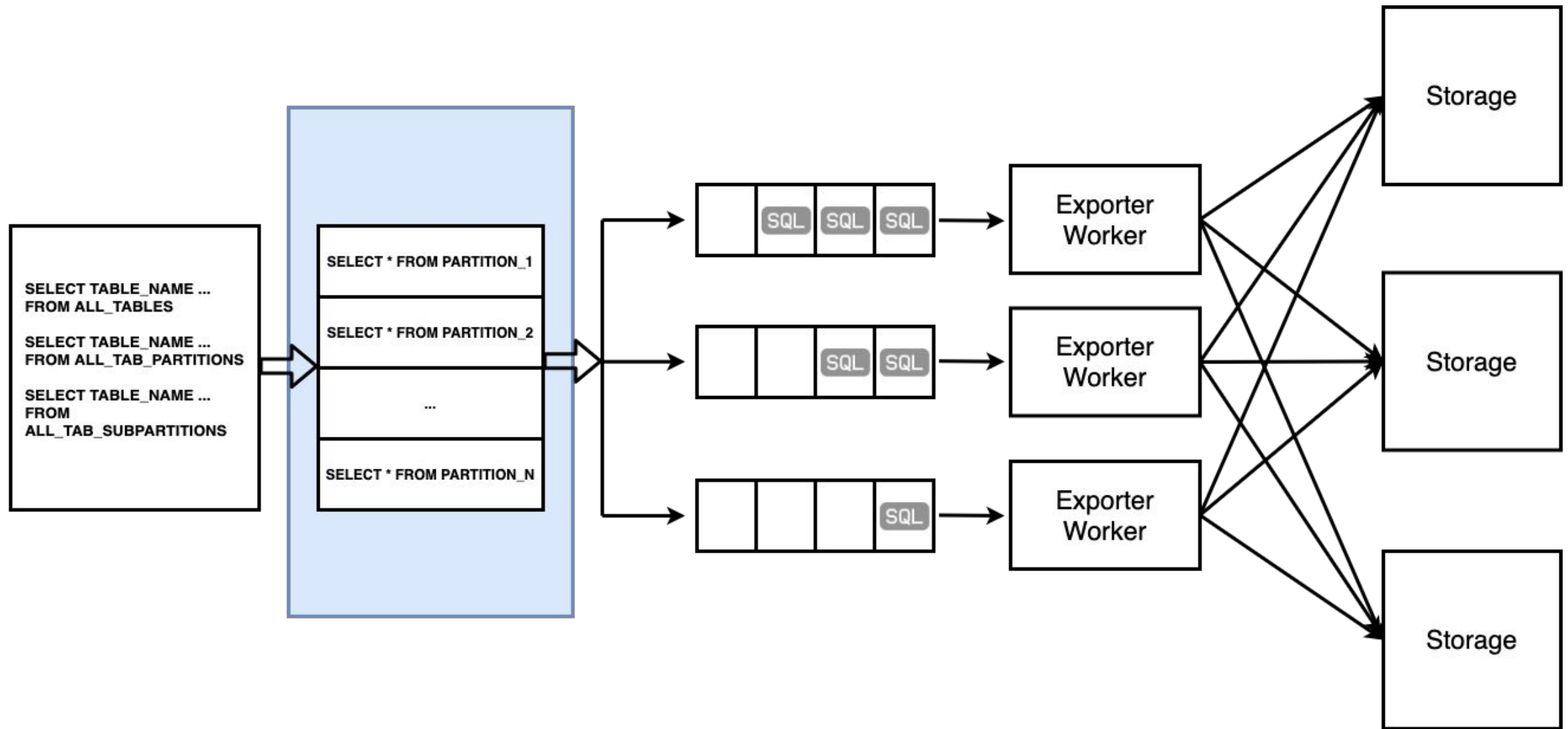


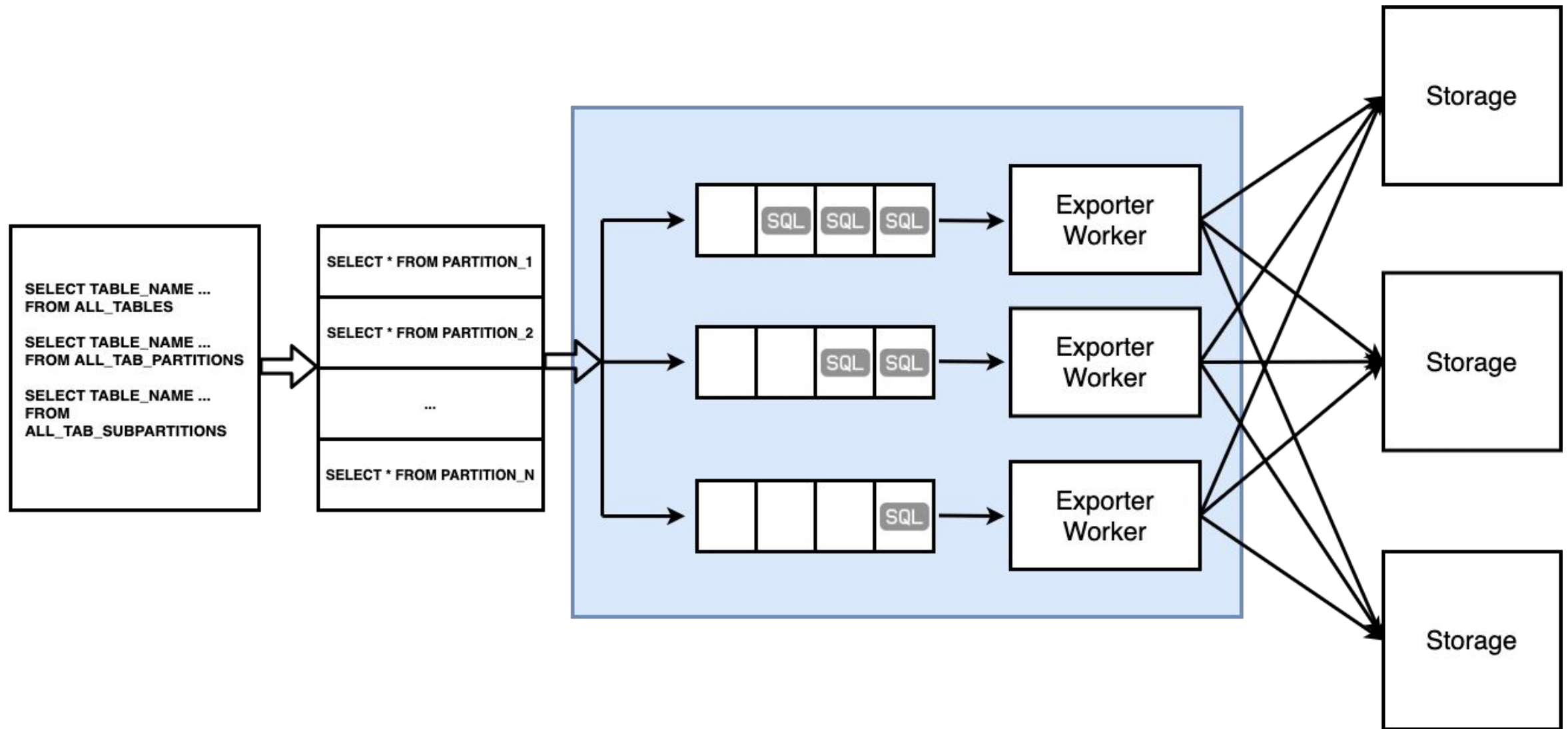
Кэши



Прогрев v2

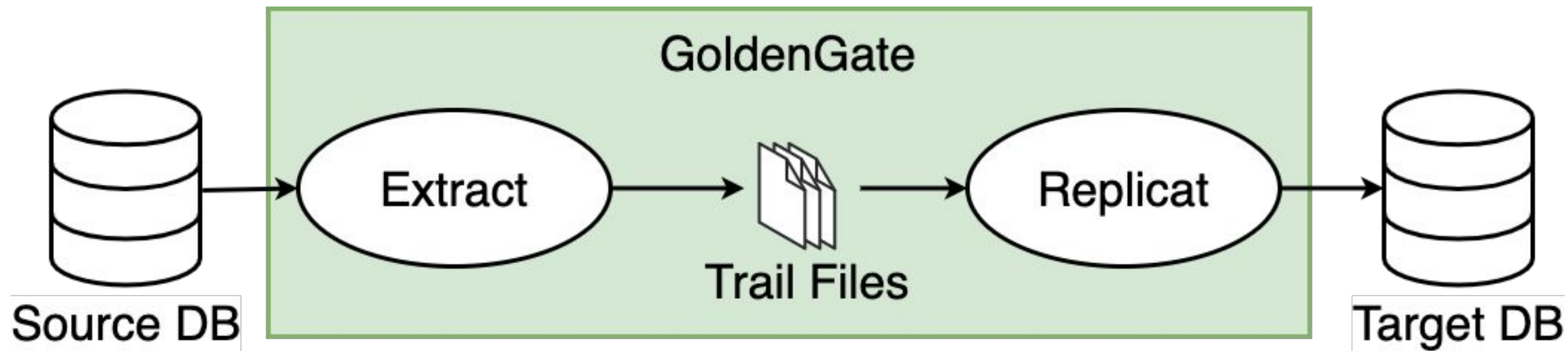


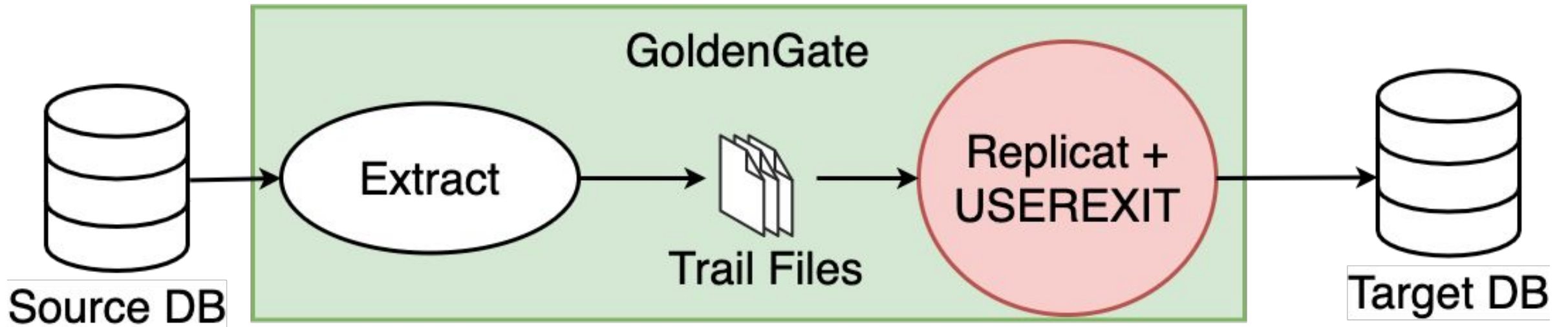


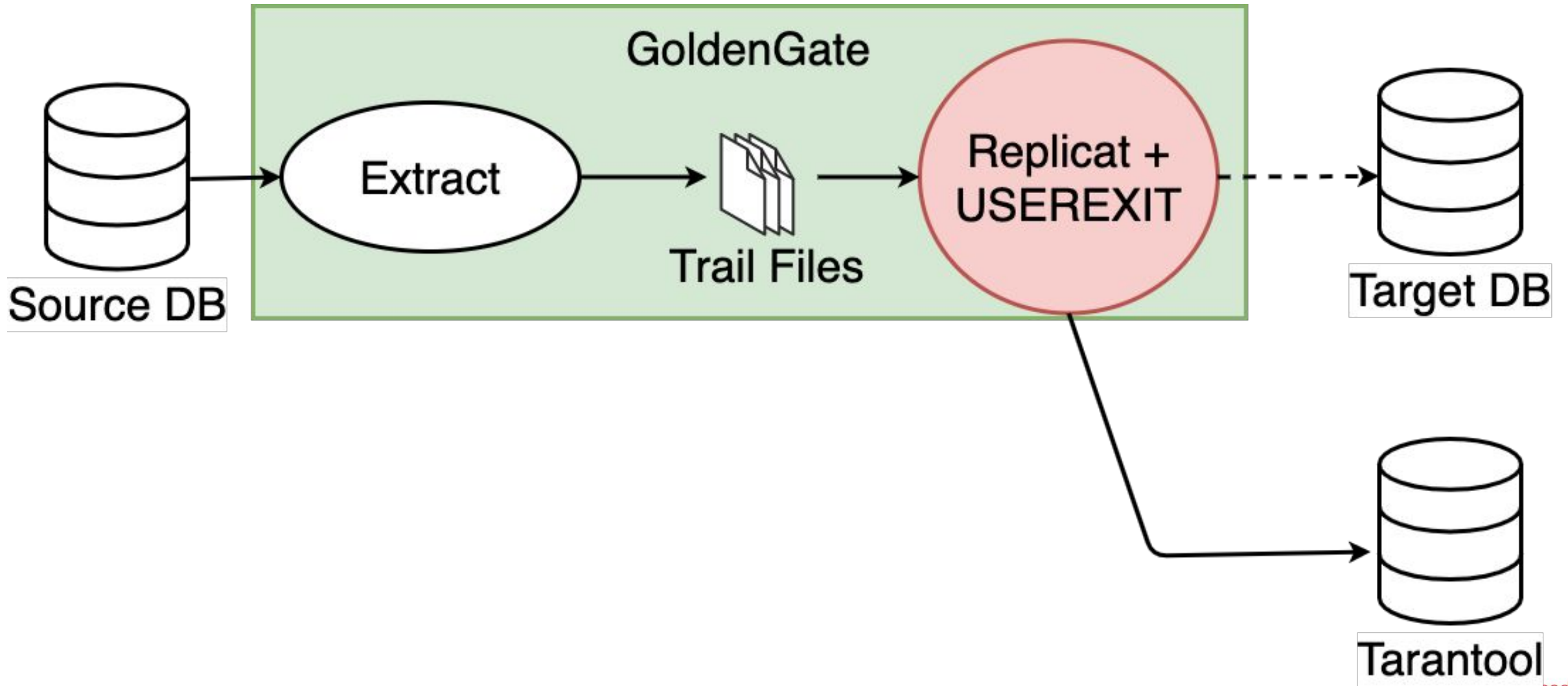


USEREXIT

Change Data Capture







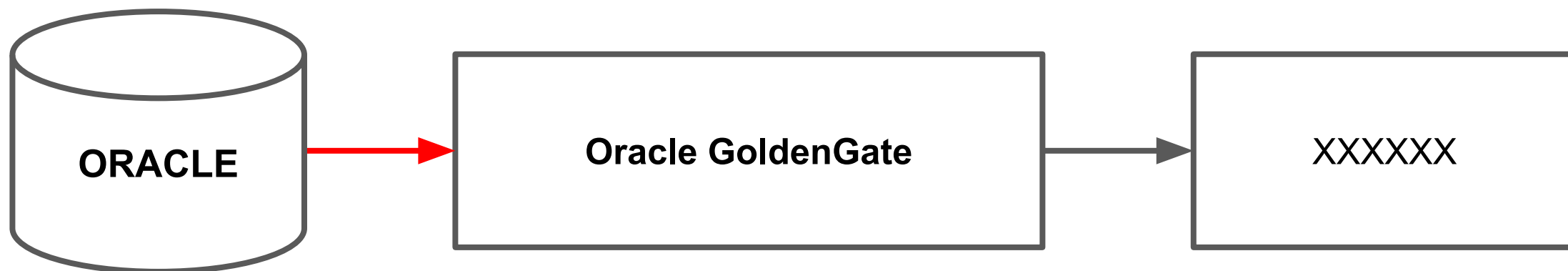
$SLA^* < 30s$

* на лаг репликаци

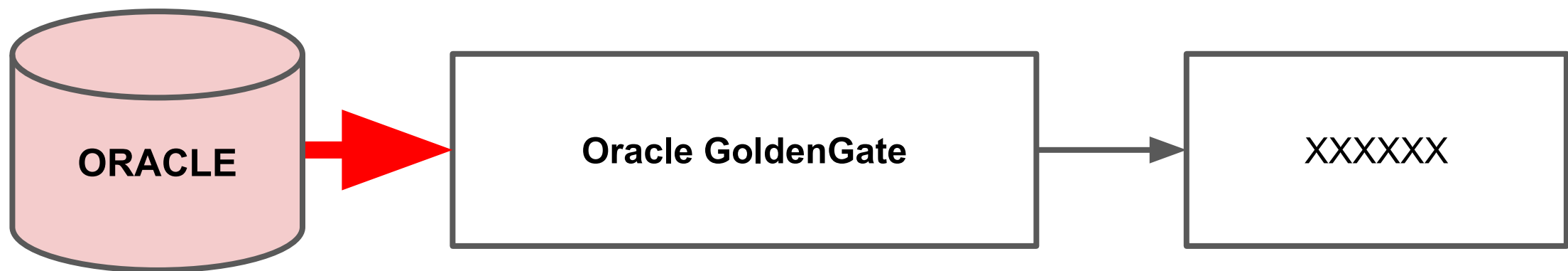
♥ Оставание репликации по времени (секунды)



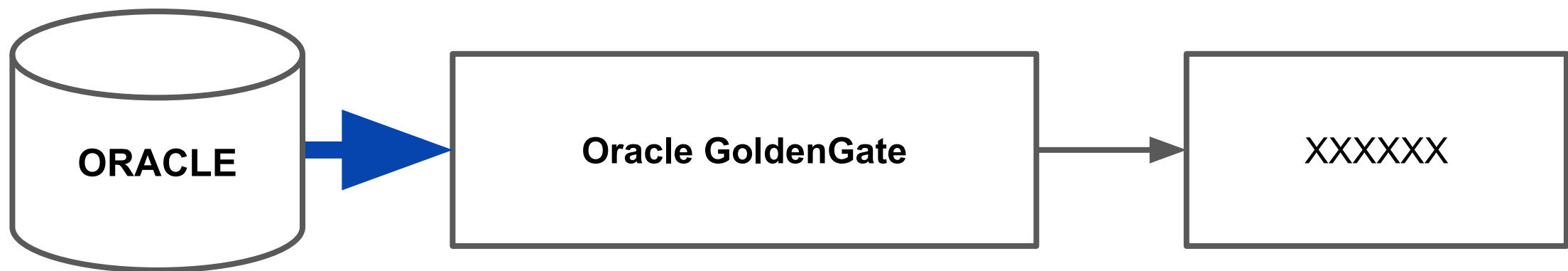
Почему так?



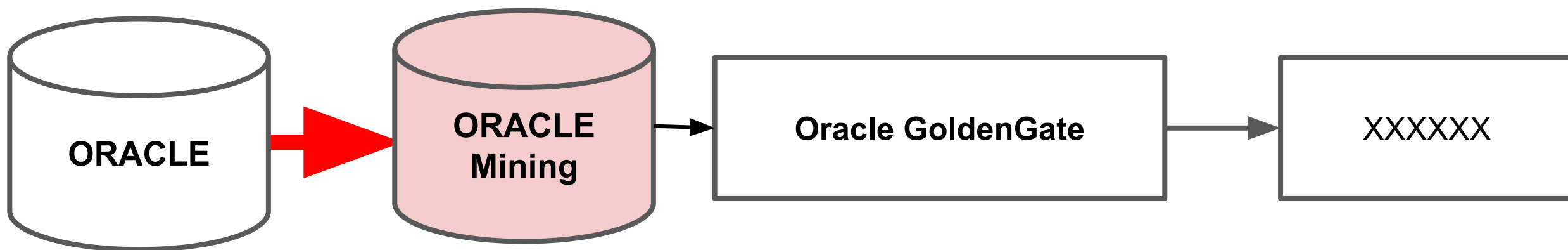
Почему так?



Почему так?



Почему так?

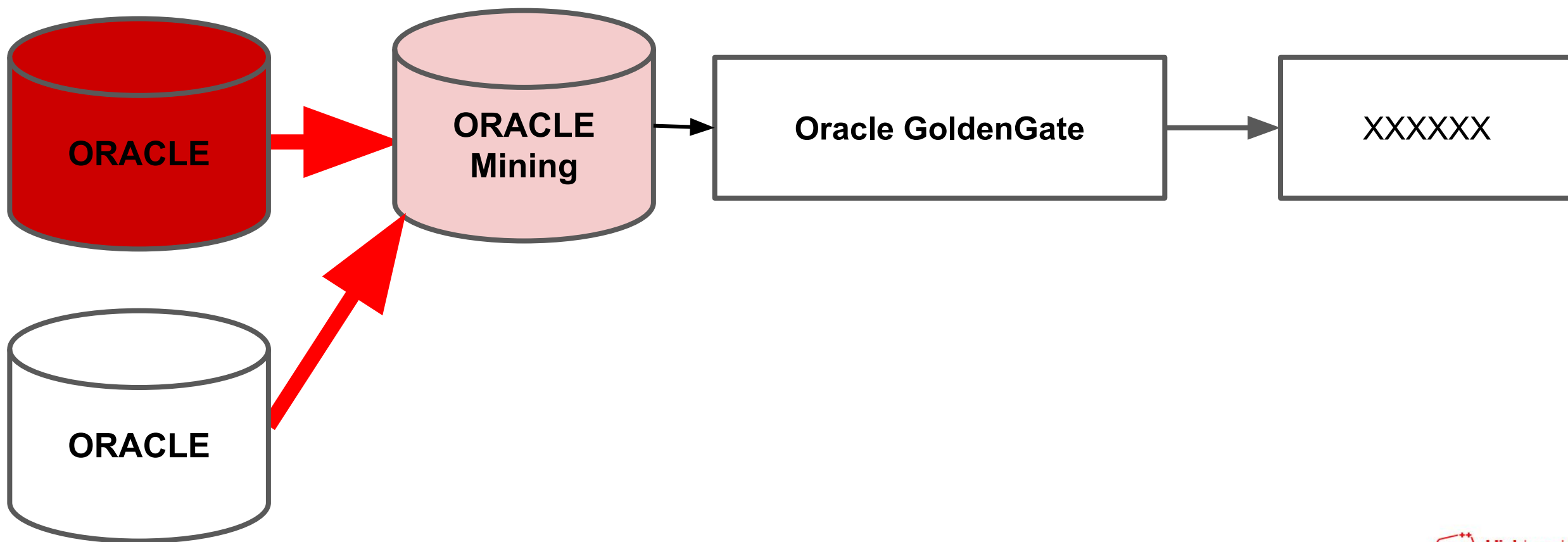


♥ Оставание репликации по времени (секунды) ▾

~18 СЕКУНД



Почему так?



Итого

- Достигли реалтайма
- Mining работает как слой изоляции от сбоев
- Нужно работать с Redo-логам, не с Archive



CP1251? UTF-8?

XML

FORMATXML ENCODING UTF-8

XML

FORMATXML ENCODING UTF-8

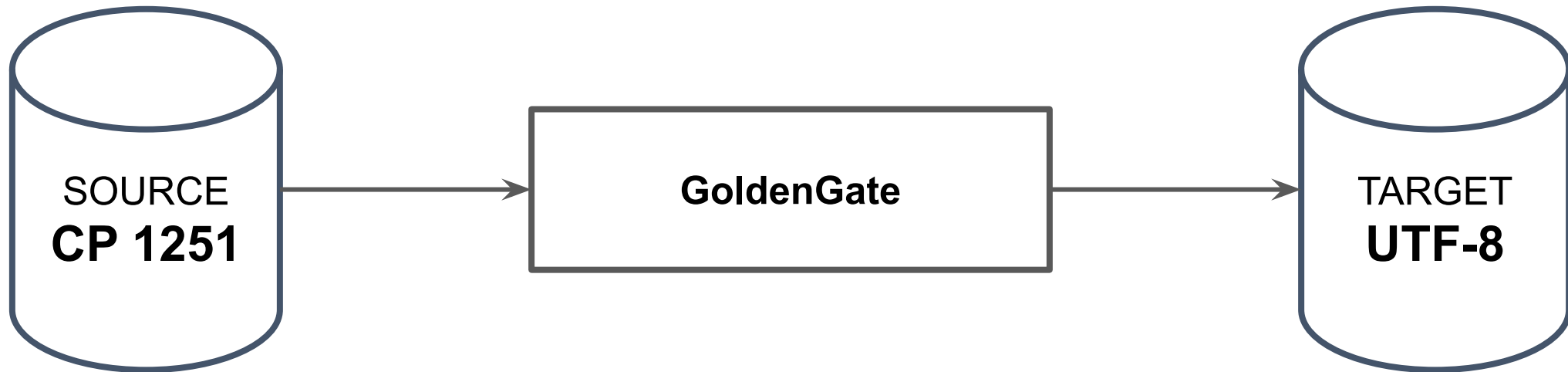
.....

<columns>

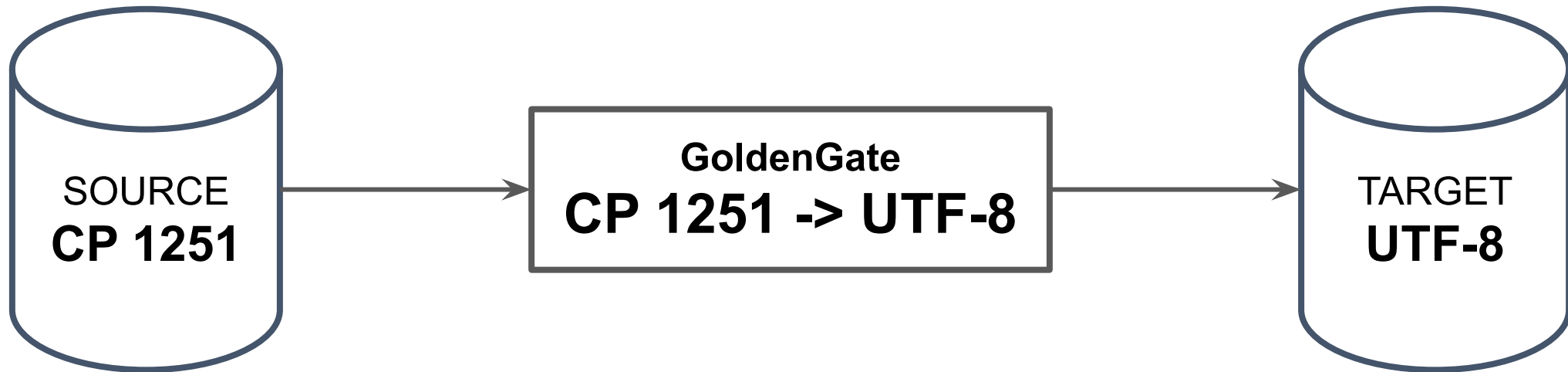
<column name="COLUMN" key=true>НОРМАЛЬНЫЙ ТАКОЙ UTF</column>

</columns>

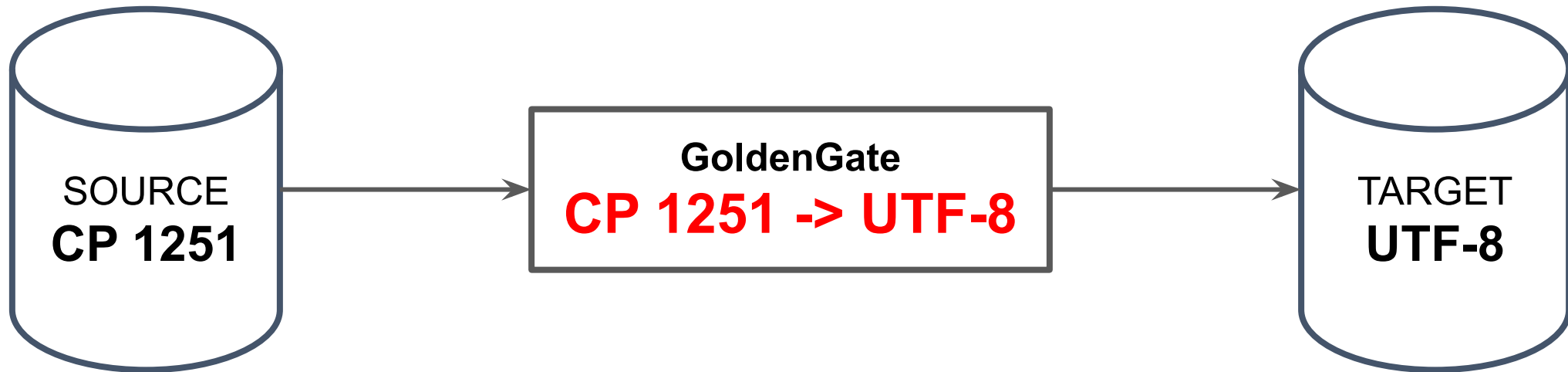
USEREXIT



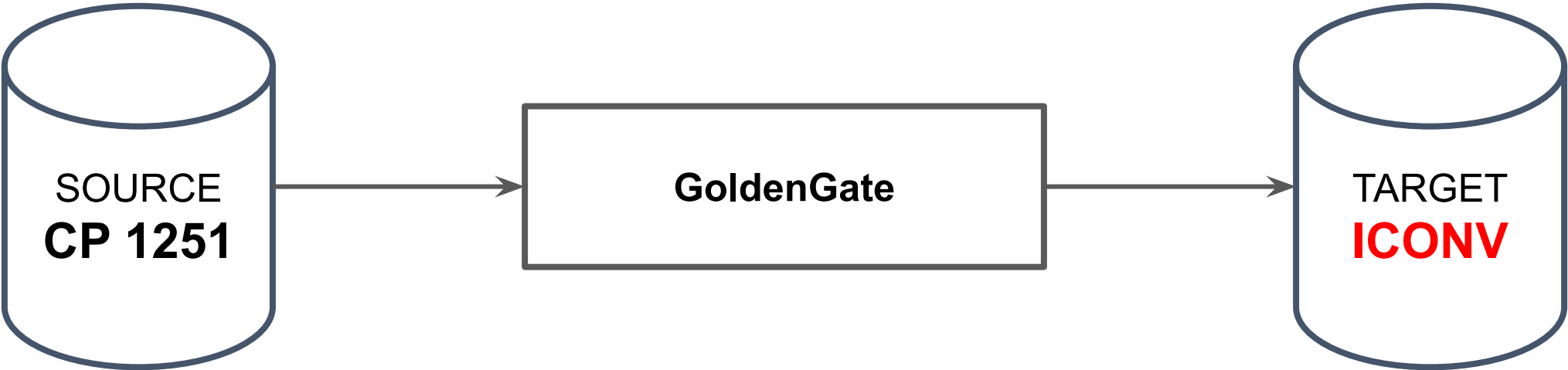
USEREXIT



USEREXIT



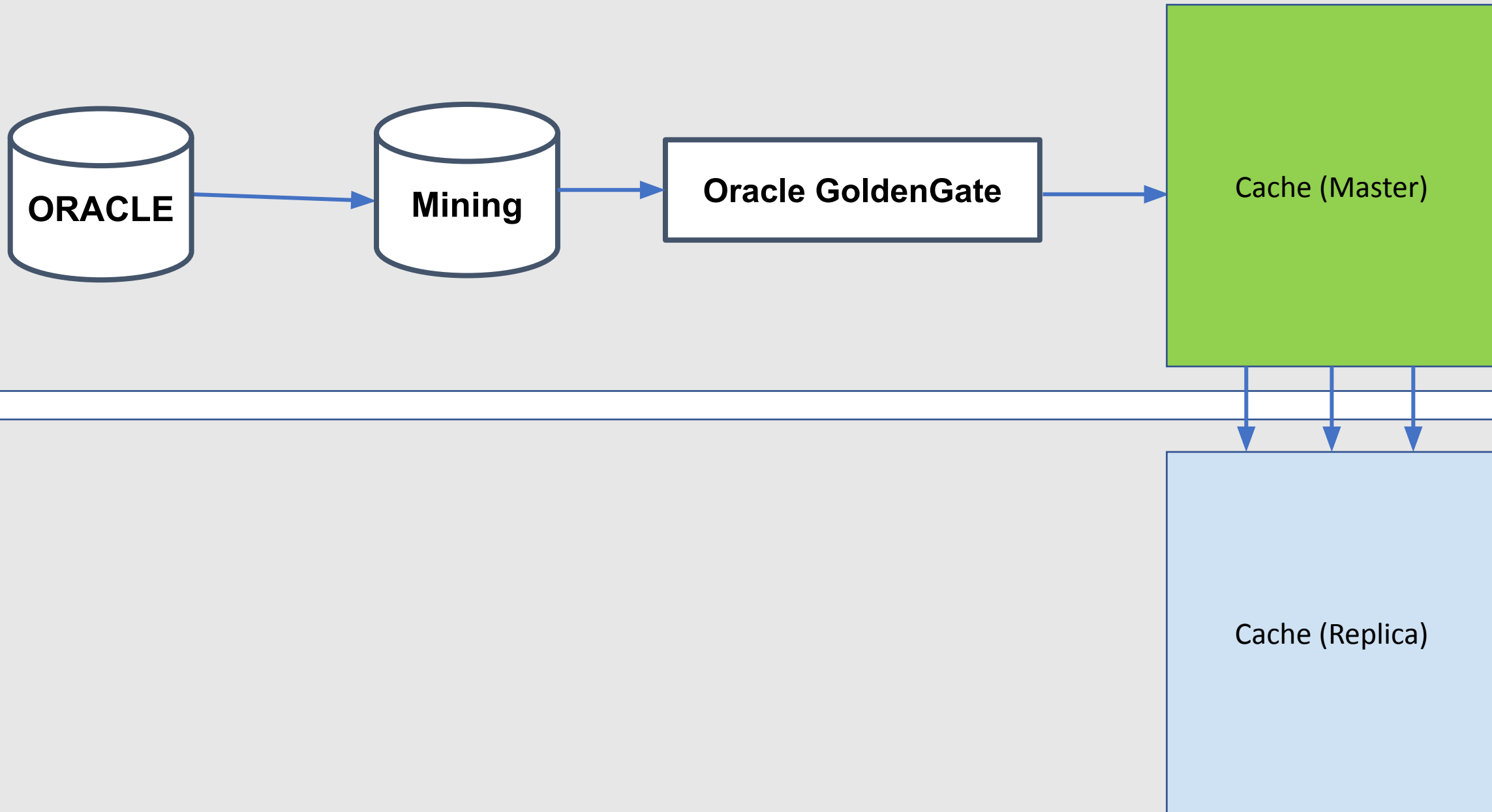
iconv

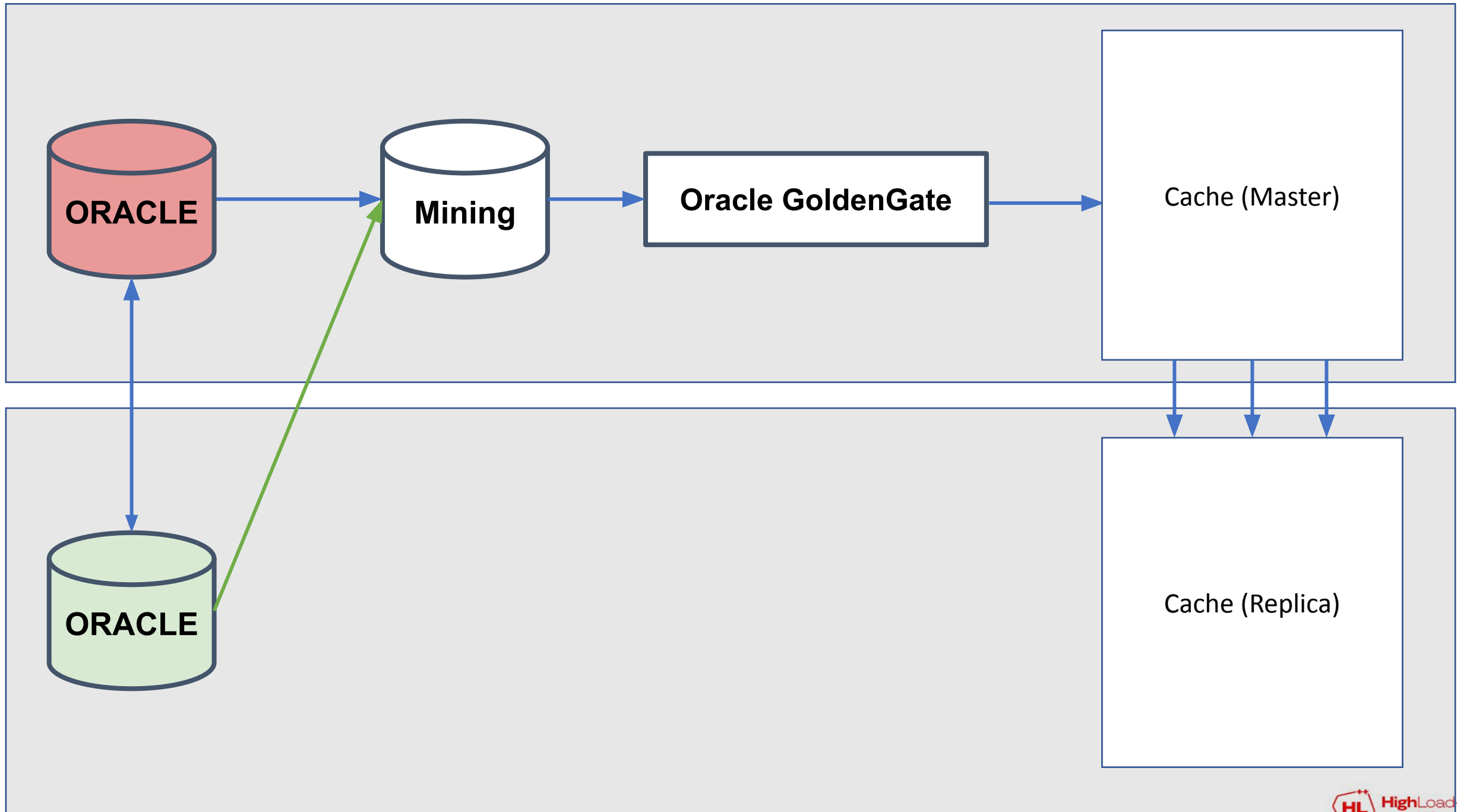


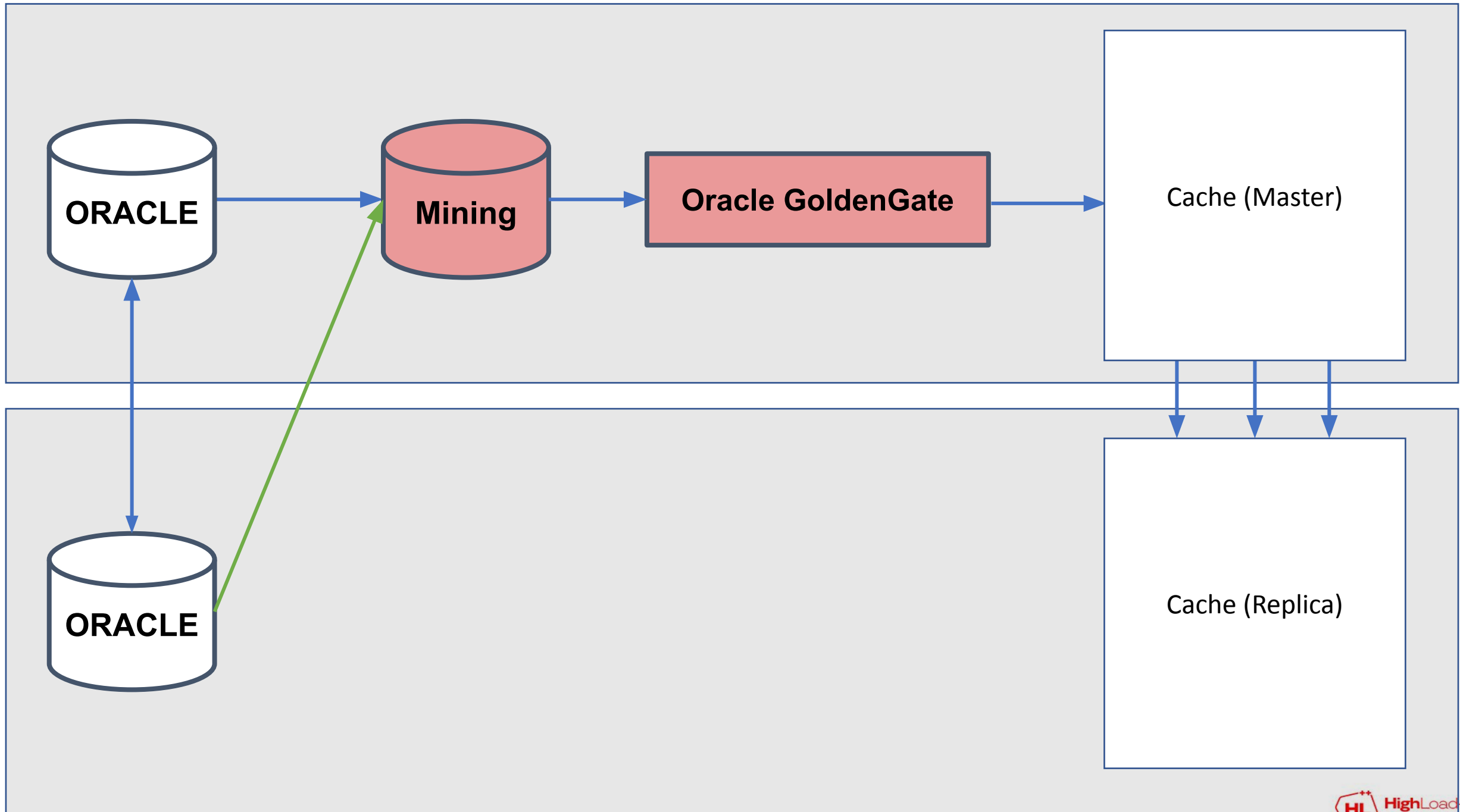
Итого

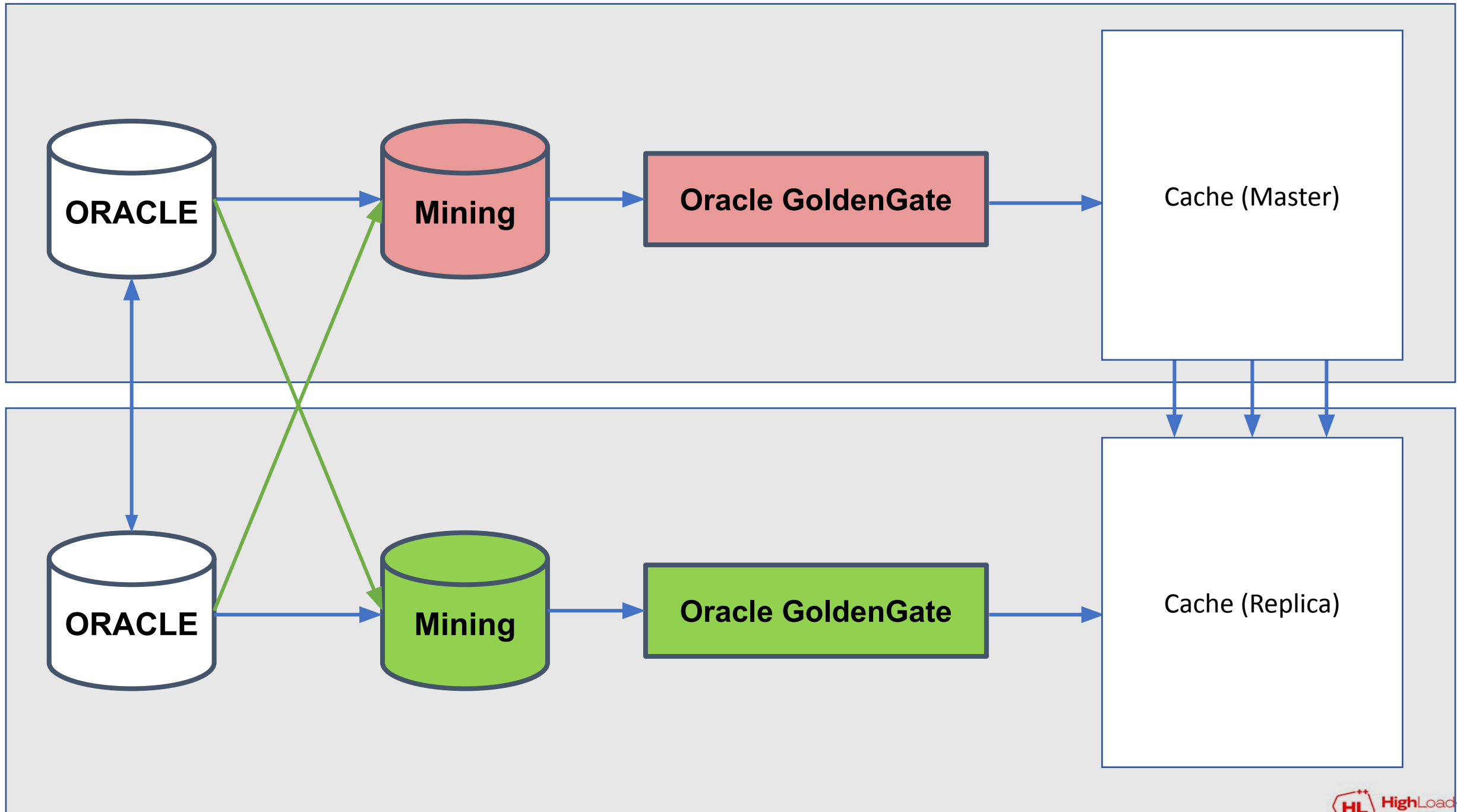
- XML — проще всего
- USEREXIT — конвертирует сам, но нам не покатило
- iconv — работает, но есть нюансы

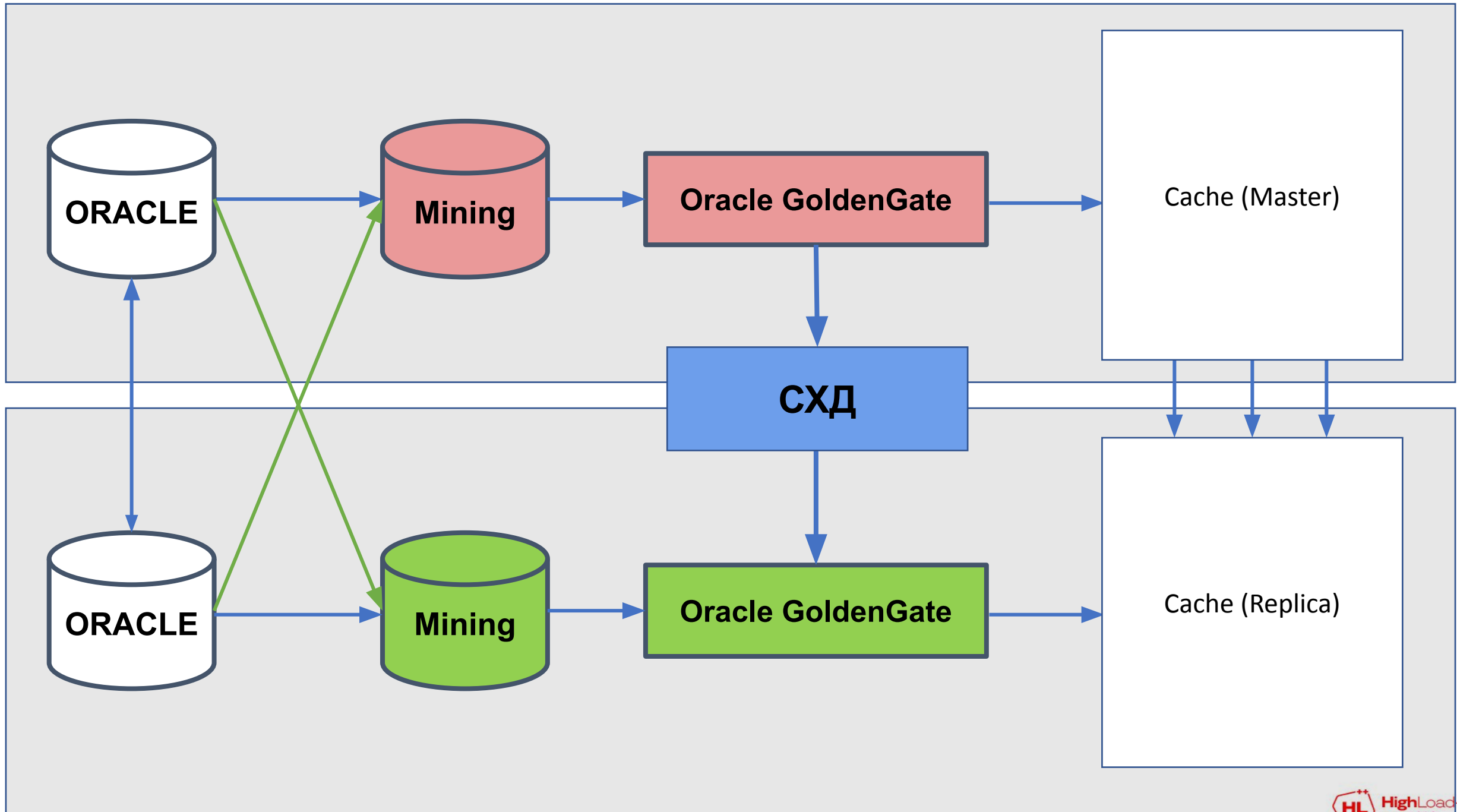
Failover

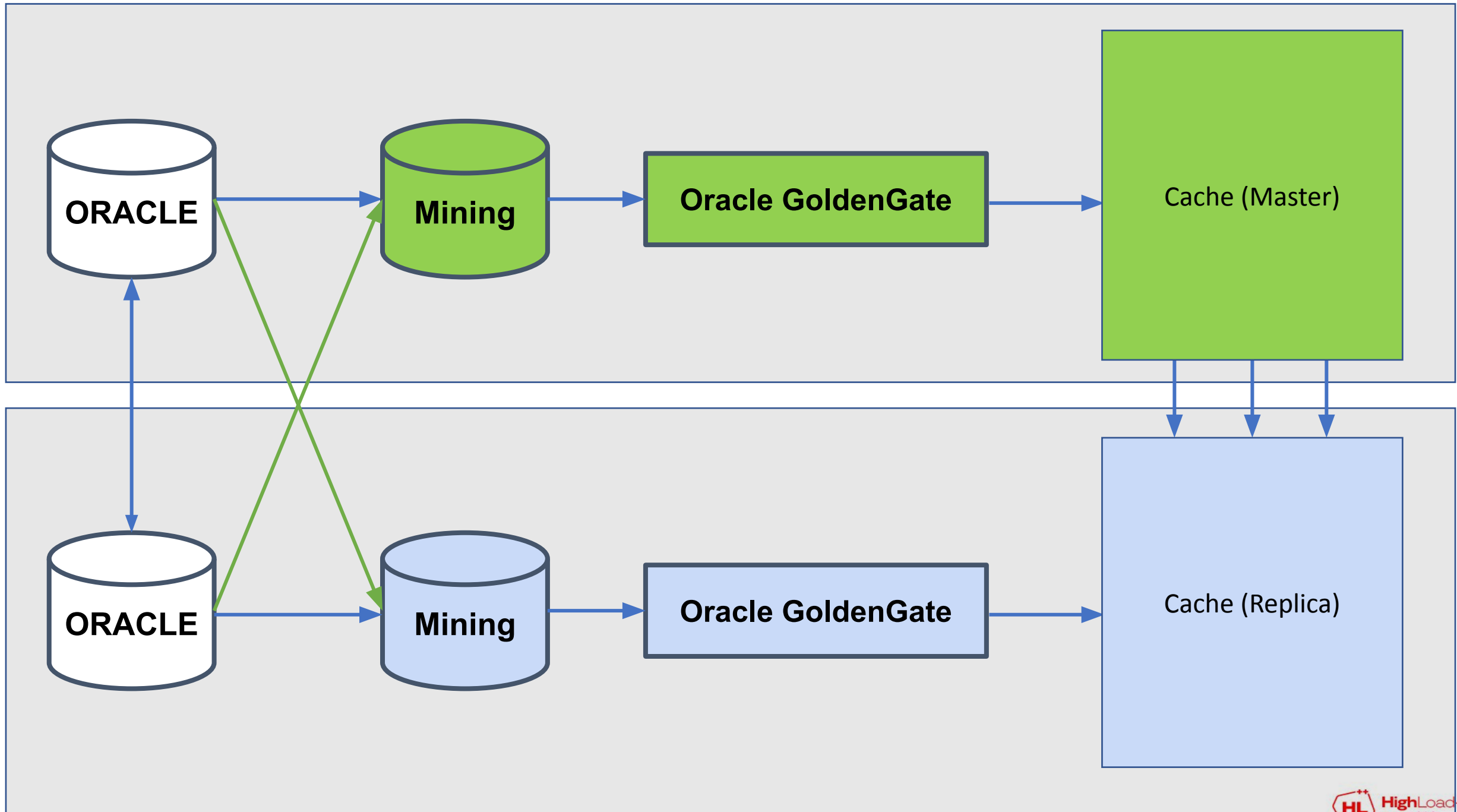


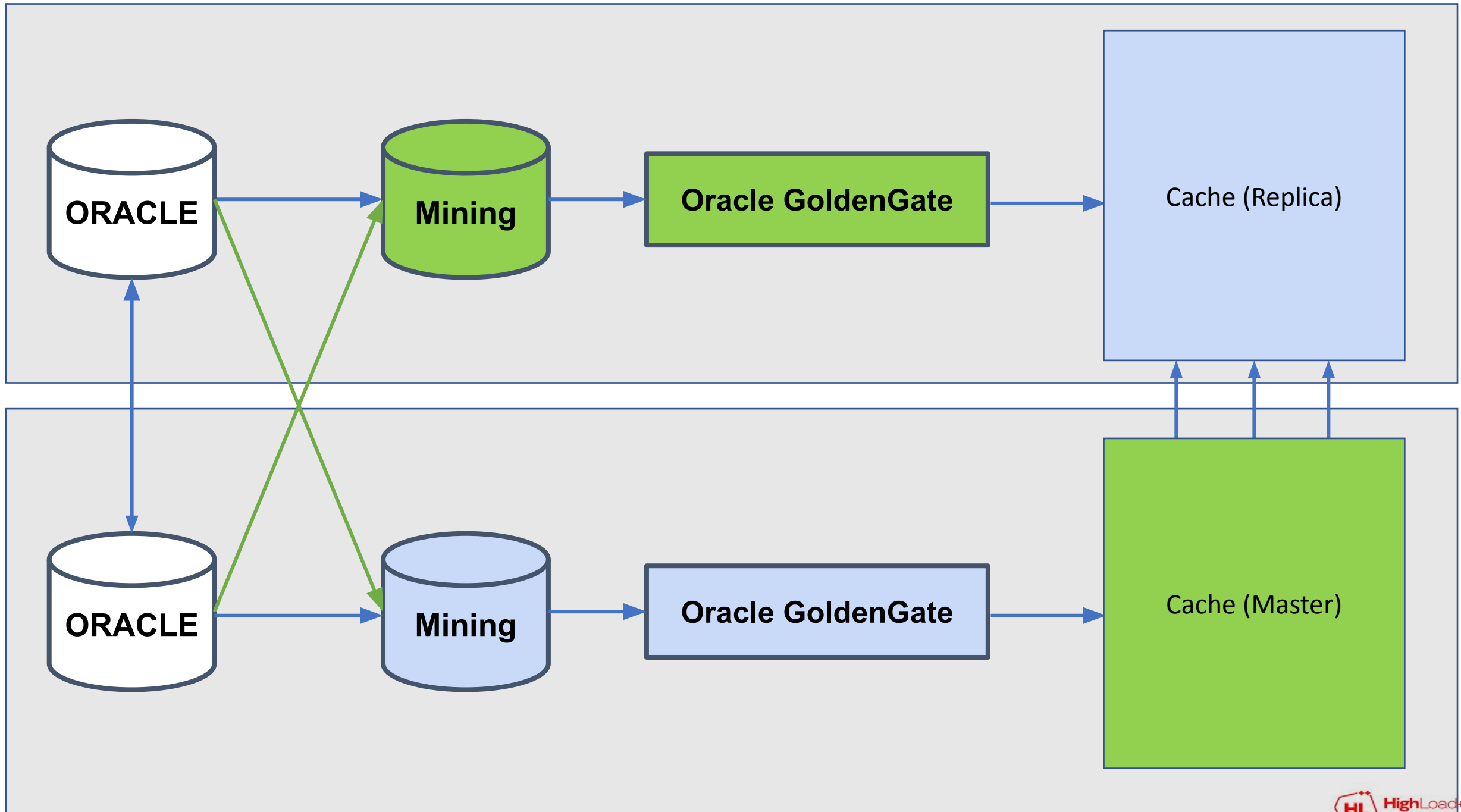


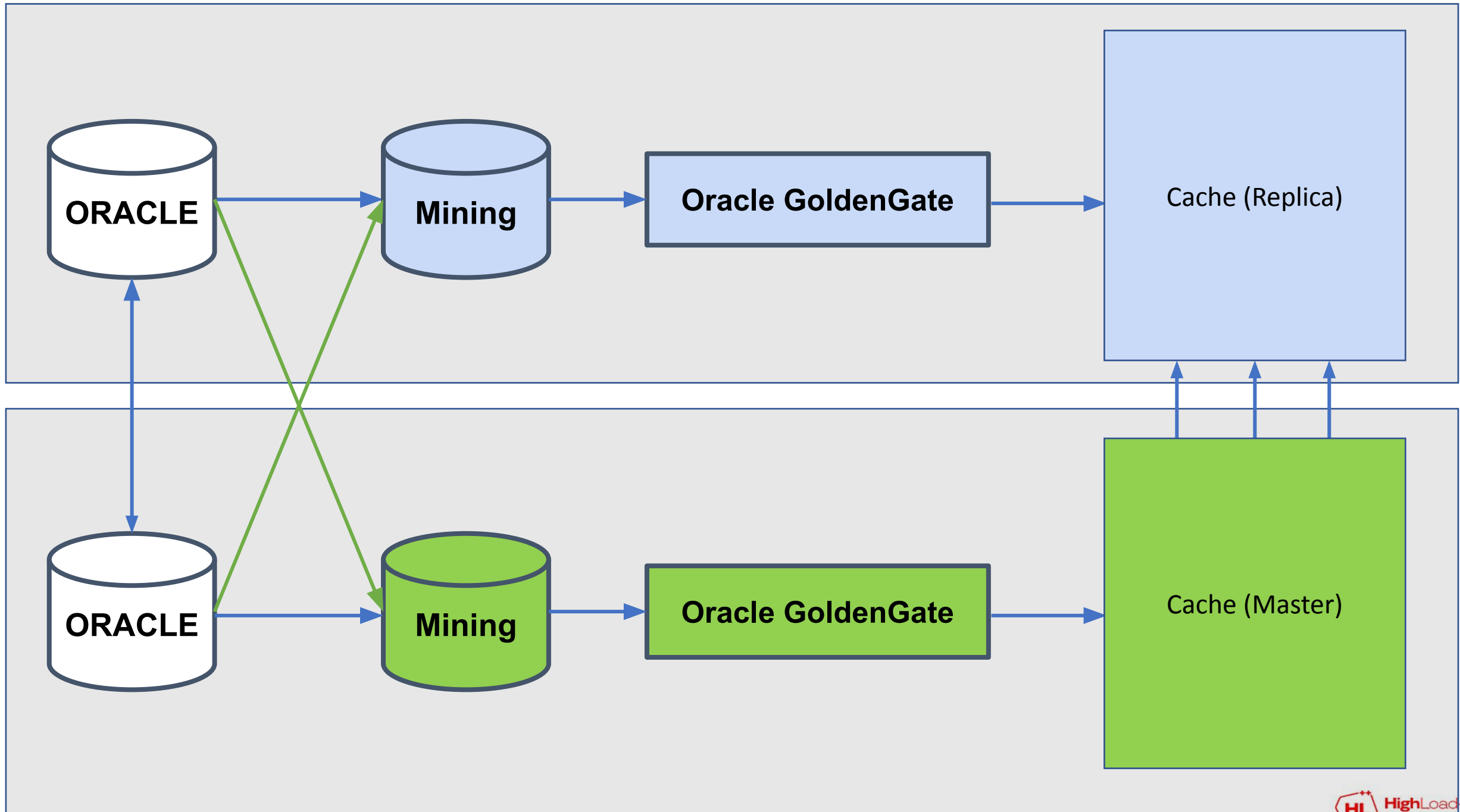










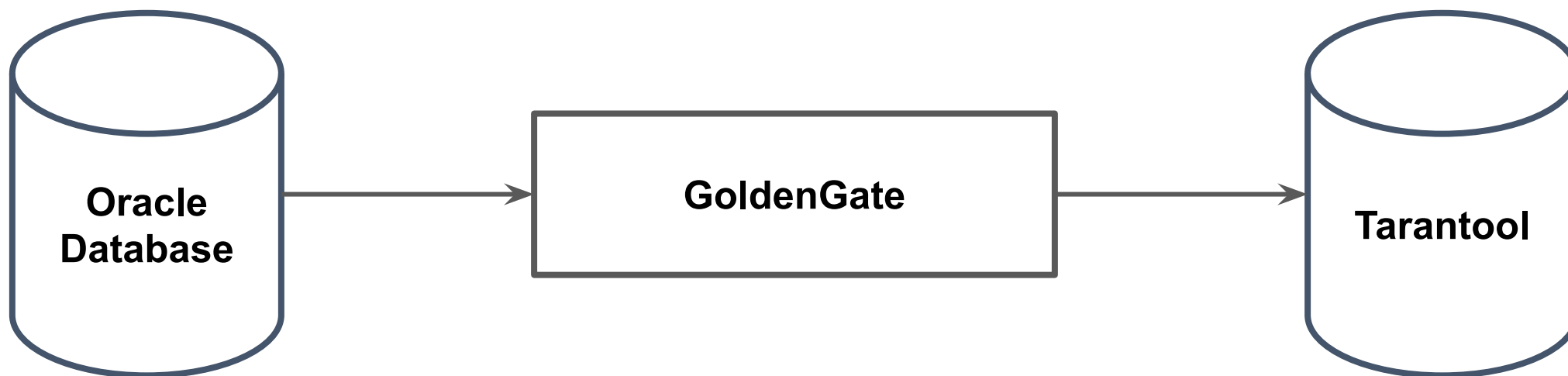


Итог

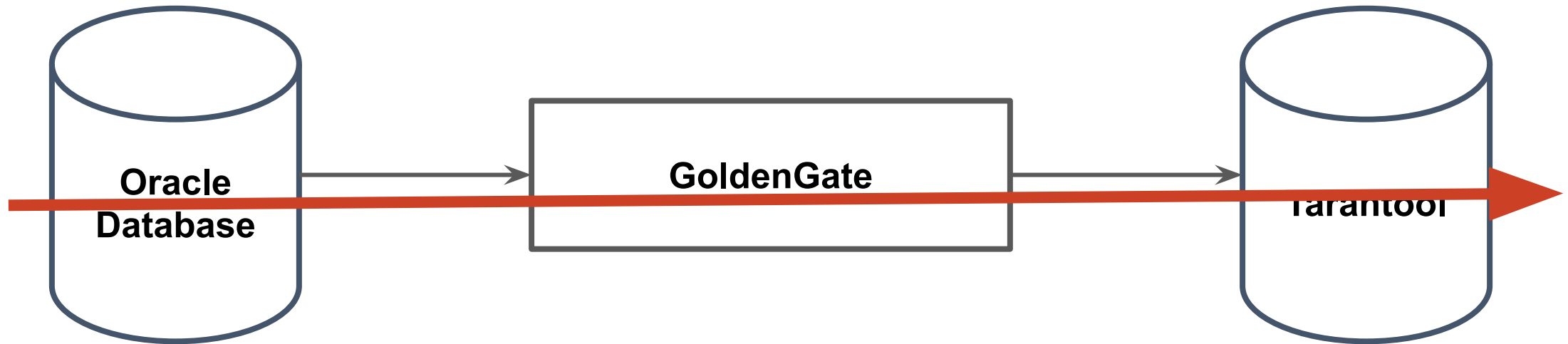
- время failover'а в рамках SLA
- фэйловер автоматический
- фэйловер GG относительно дешевый

Немного про эксплуатацию

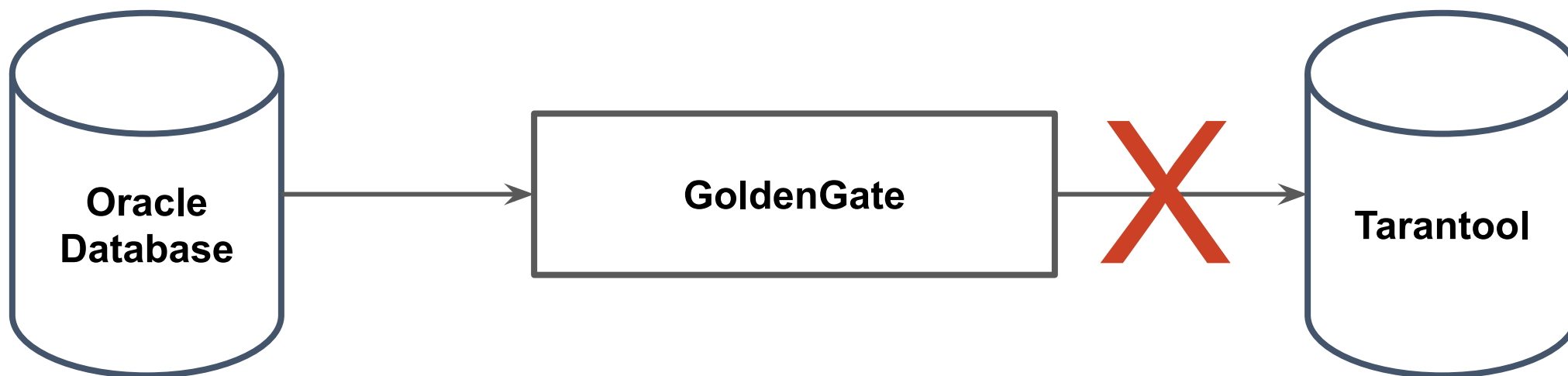
Нет “обратной связи”



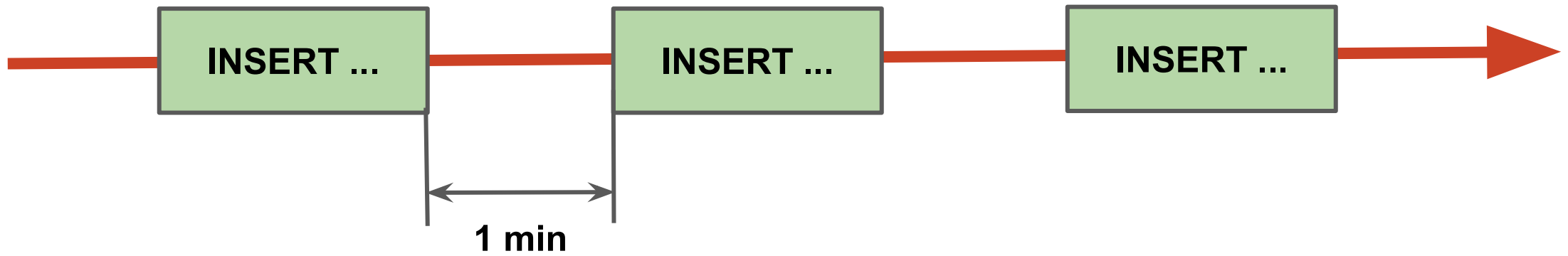
Нет “обратной связи”



Нет “обратной связи”



Heartbeat



Мониторинг: лаг репликации end 2 end

<tokens>

<token name="X_COMMIT_TIME">...</token>

</tokens>

LAG = NOW - X_COMMIT_TIME

Итог

- **Непрерывность репликации**
- **Лег репликации**
- Много чего еще — статусы всех процессов, нагрузка на downstream базе, etc.

Итого по технике

- Характеристики источника
- Структура витрины
- XML — легальный и самый простой способ CDC
- USEREXIT — более технологичный способ CDC
- Важны эксплуатационные характеристики

Итог



Tarantool на Highload